

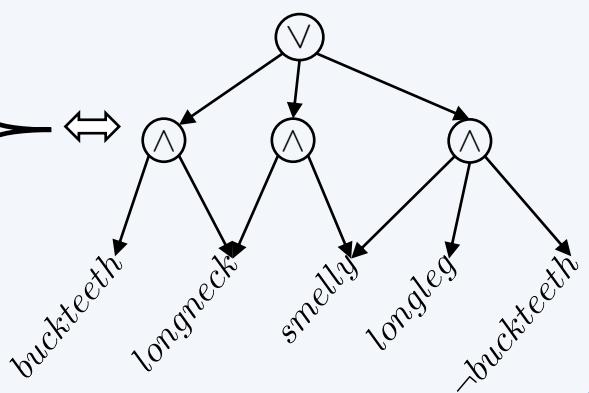
Class: Horse

MoIE

smelly ↔

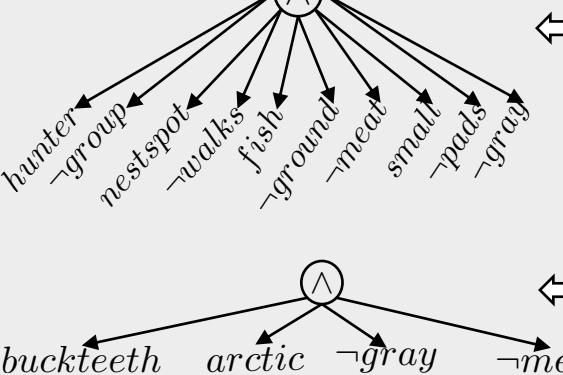


Baseline



Class: Otter

MoIE



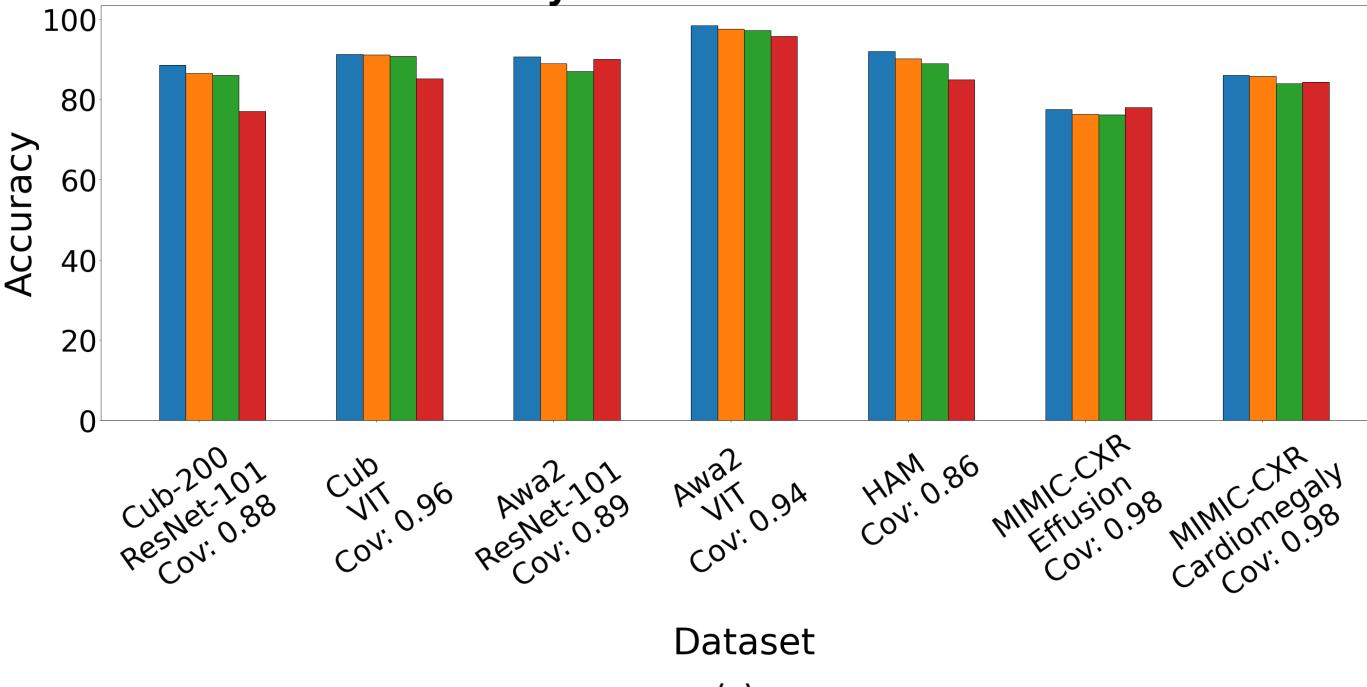
Baseline



(a)

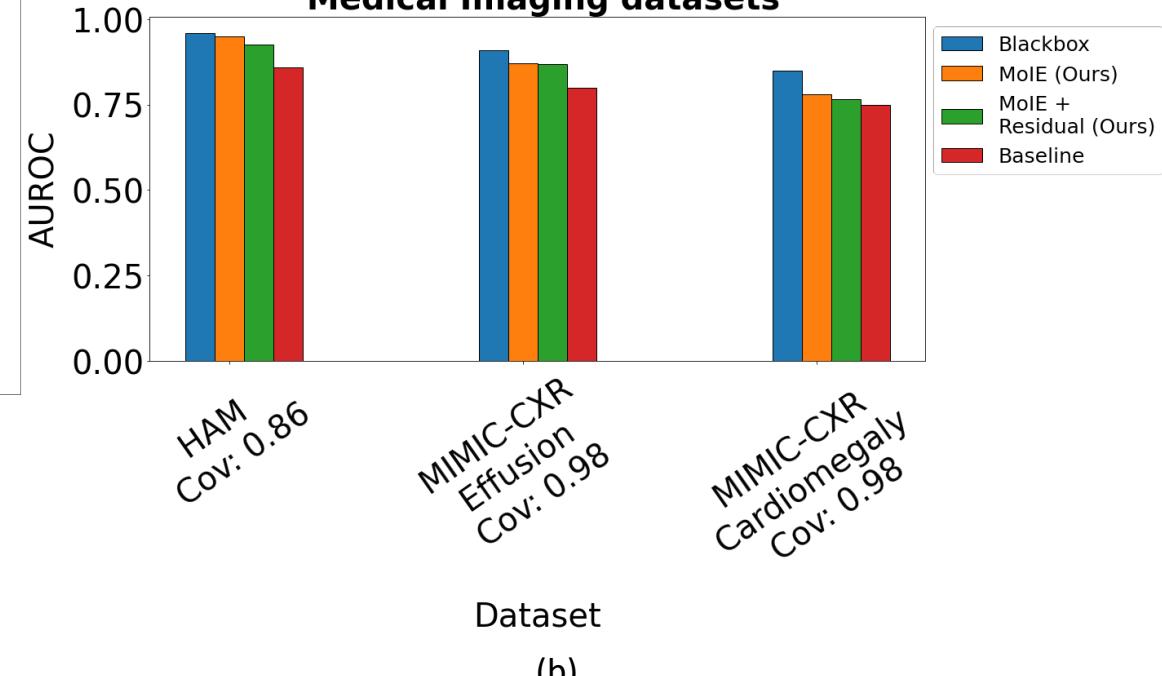
(b)

Accuracy scores of all the datasets



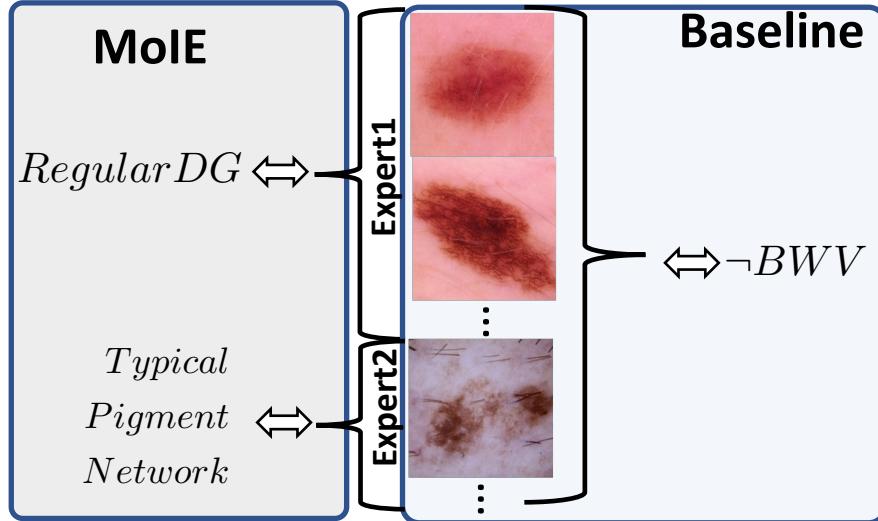
(a)

AUROC scores of the Medical Imaging datasets



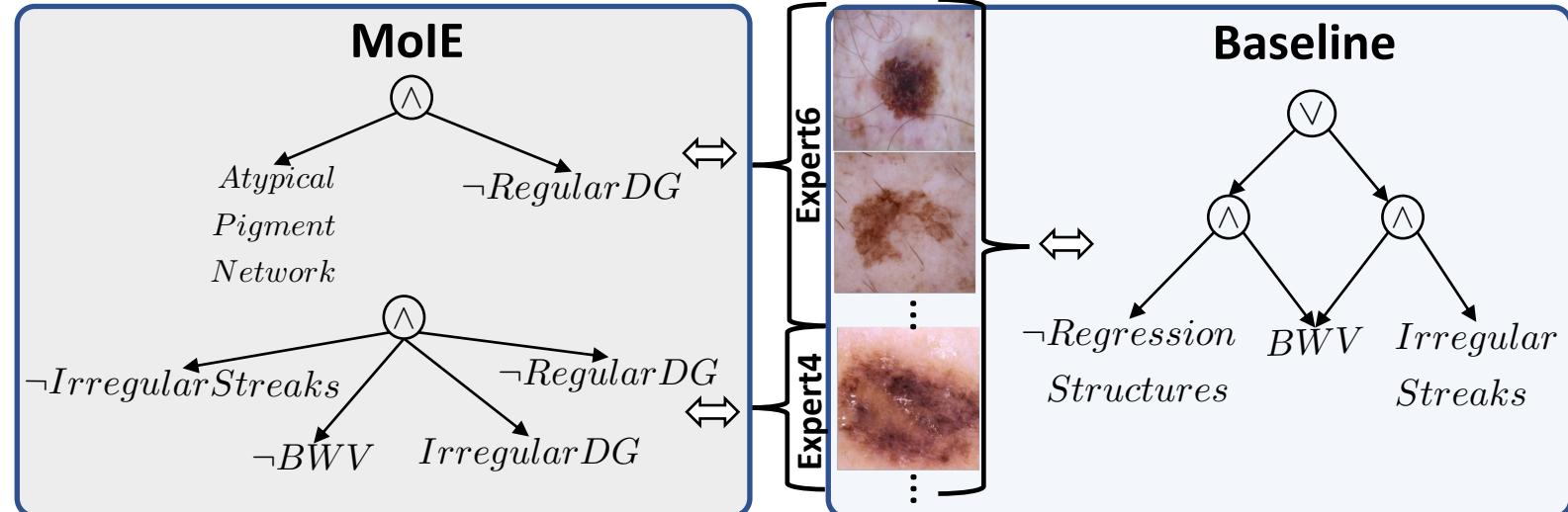
(b)

Class: Benign

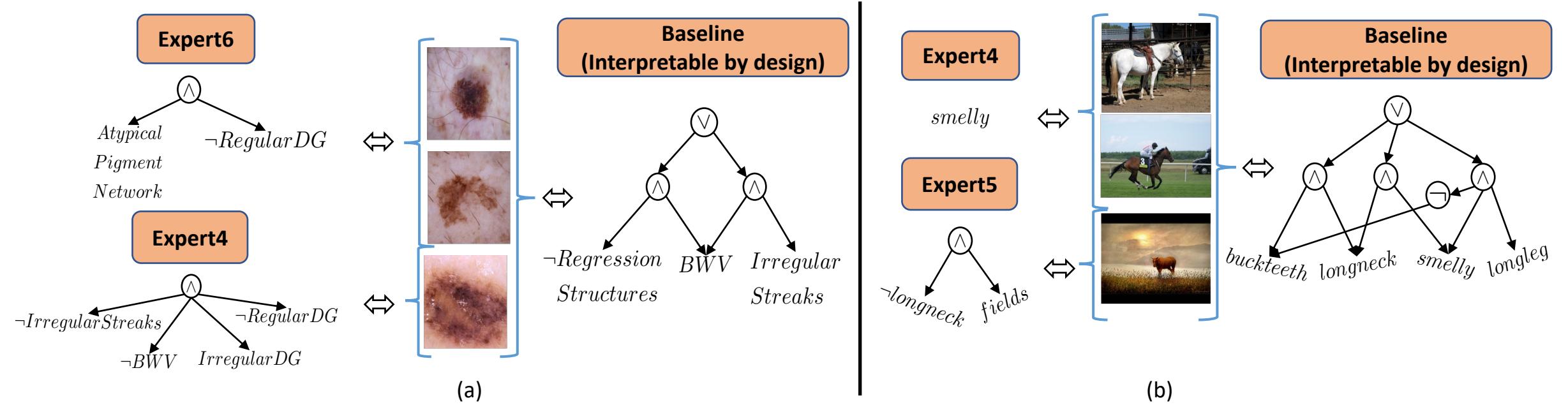


(c)

Class: Malignant



(d)



Class: Horse

MoIE

smelly ↔

Expert4



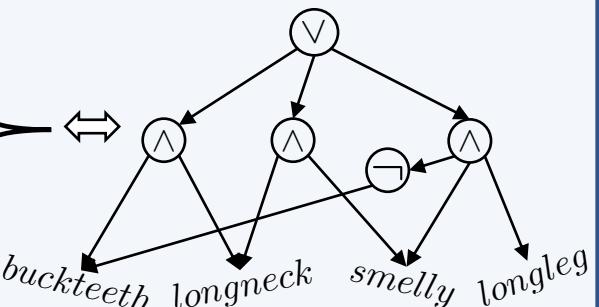
⋮

Expert5



⋮

Baseline



Class: Otter

MoIE

hunter
¬group
nestspot
¬walks
fish
¬ground
¬meat
small
¬pads
¬gray

buckteeth
arctic
¬gray
¬meat

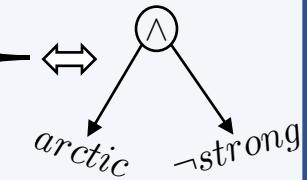
Expert1



⋮

Expert2

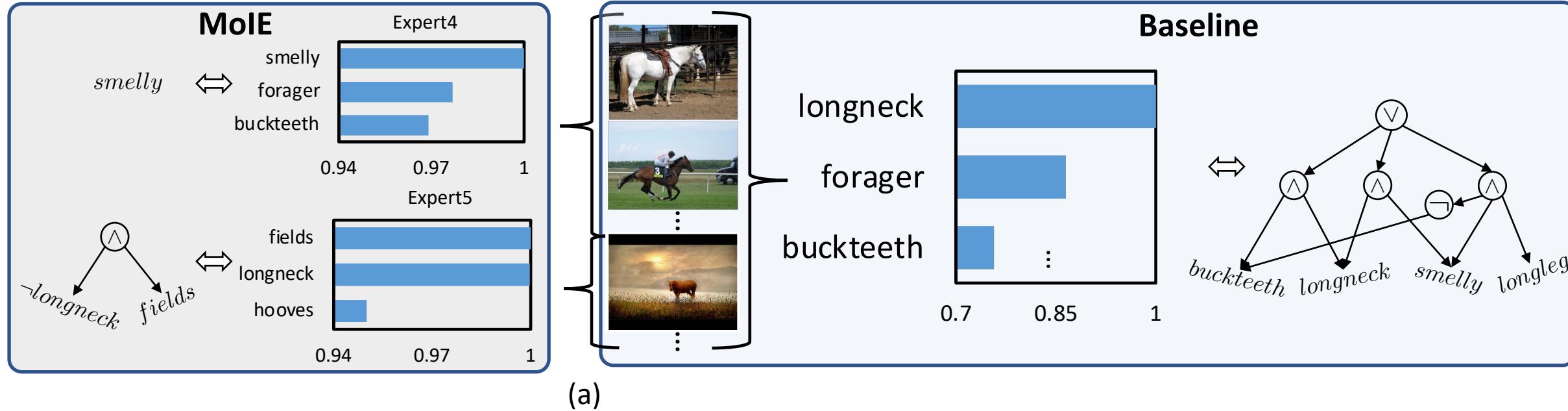
Baseline

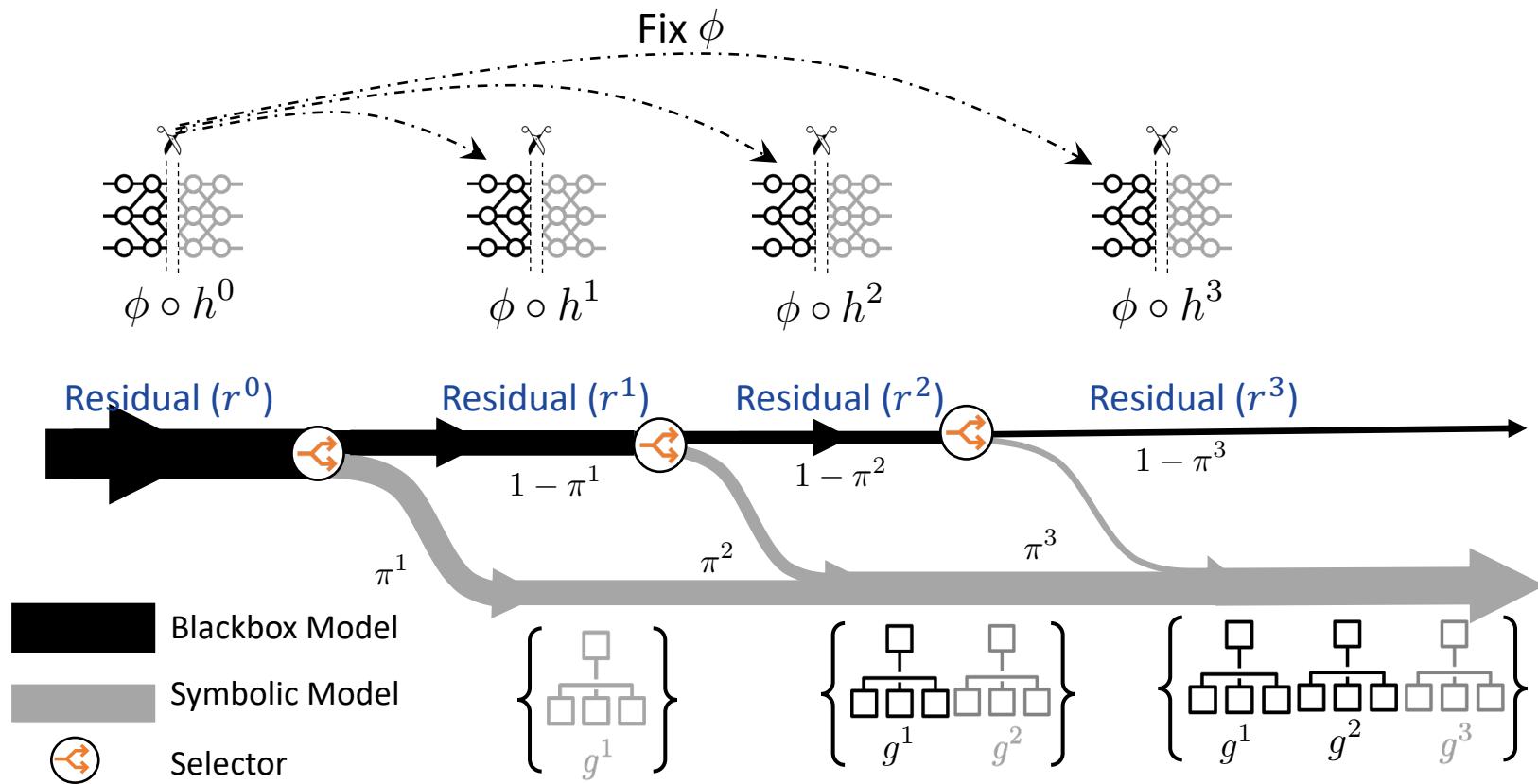


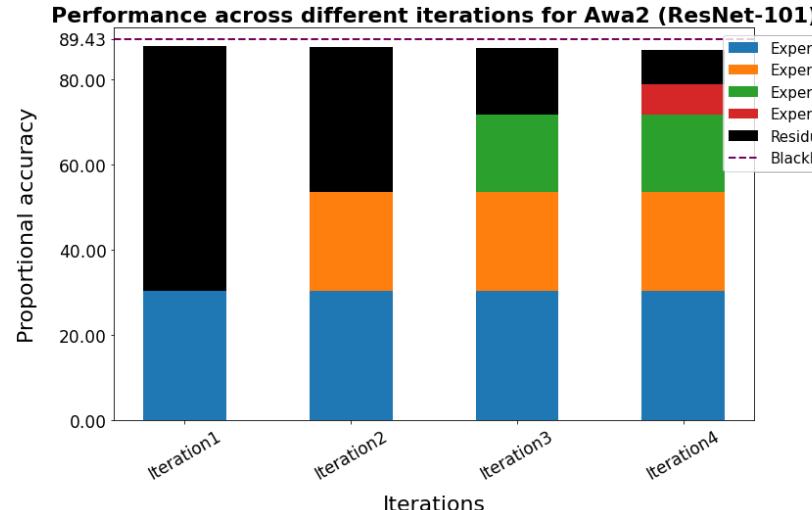
(a)

(b)

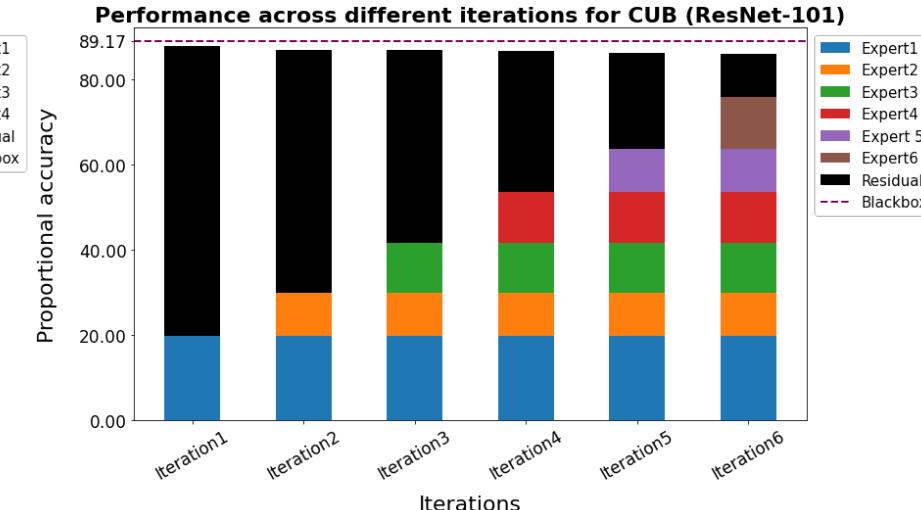
Class: Horse



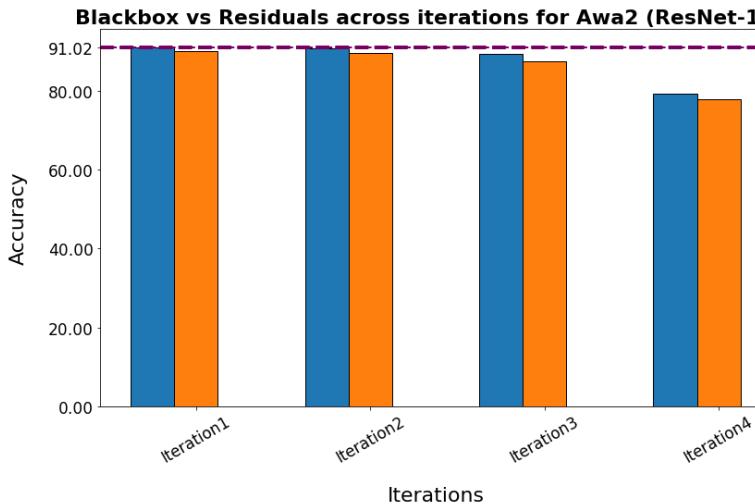




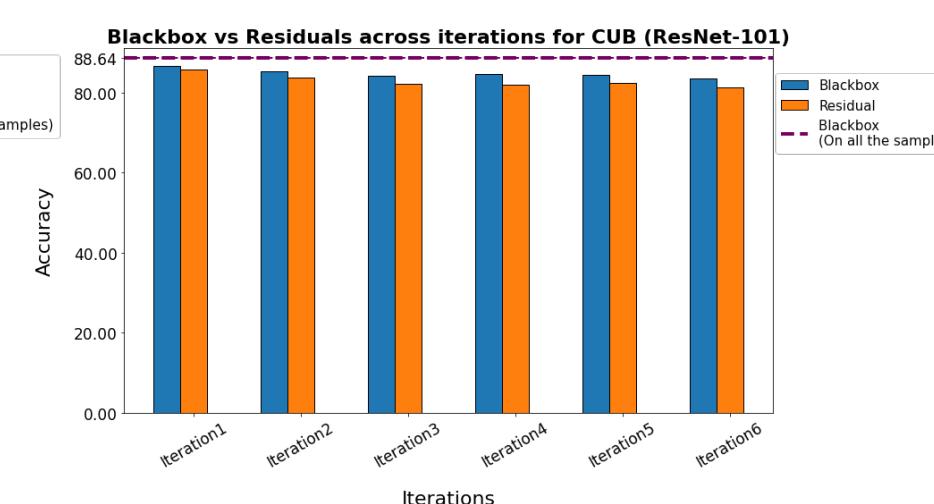
(a)



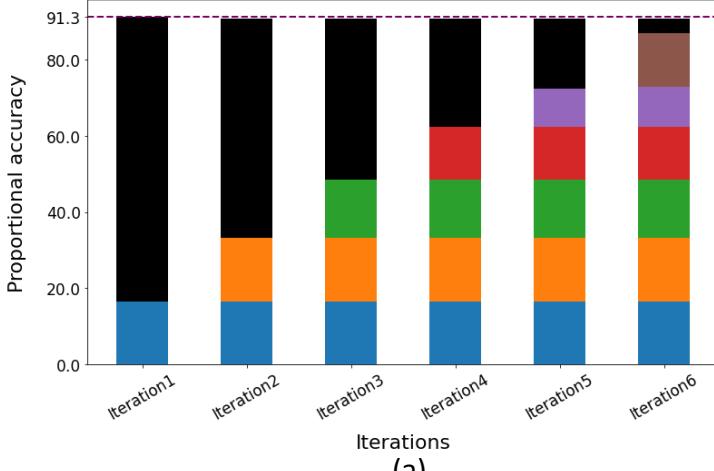
(b)



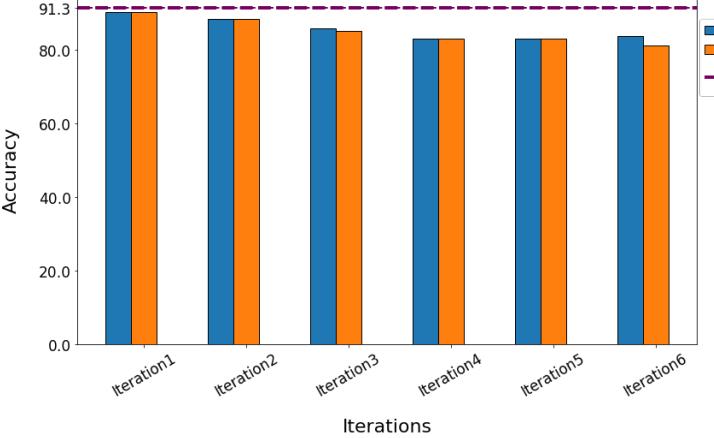
(c)



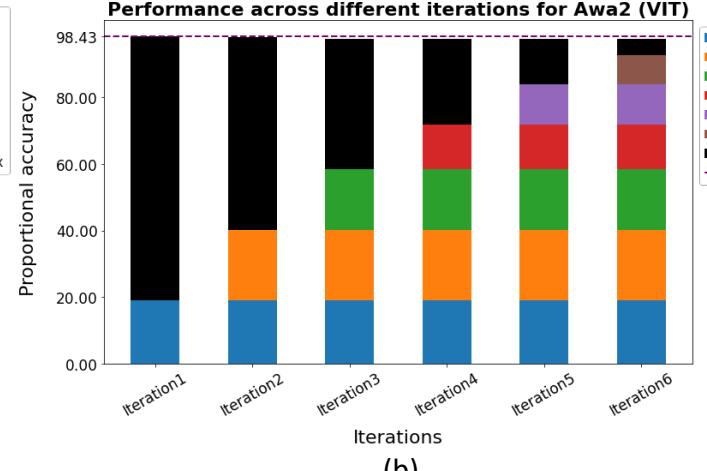
(d)

Coverages of various experts for CUB (ViT)**Coverages of various experts for Awa2 (ViT)****Coverages of various experts for HAM****Performance across different iterations for CUB (ViT)**

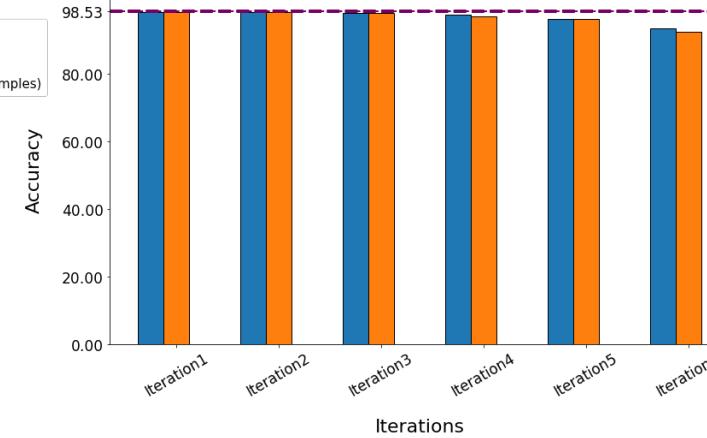
(a)

Blackbox vs Residuals across iterations for CUB (ViT)

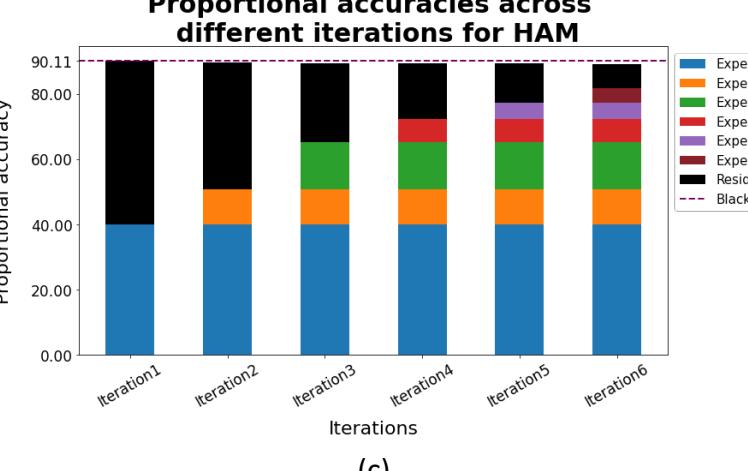
(d)

Performance across different iterations for Awa2 (ViT)

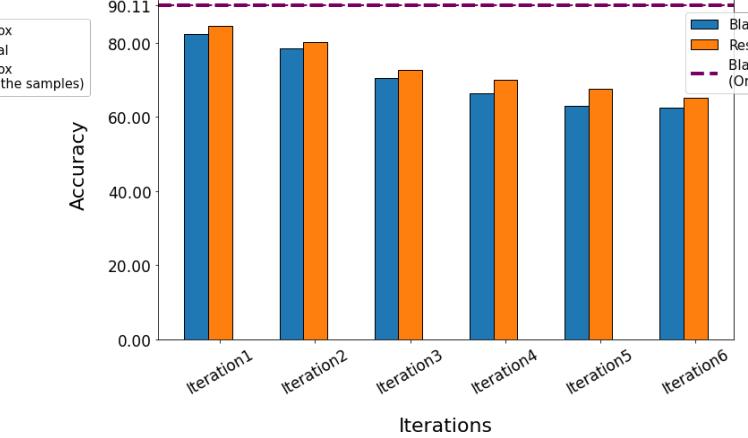
(b)

Blackbox vs Residuals across iterations for Awa2 (ViT)

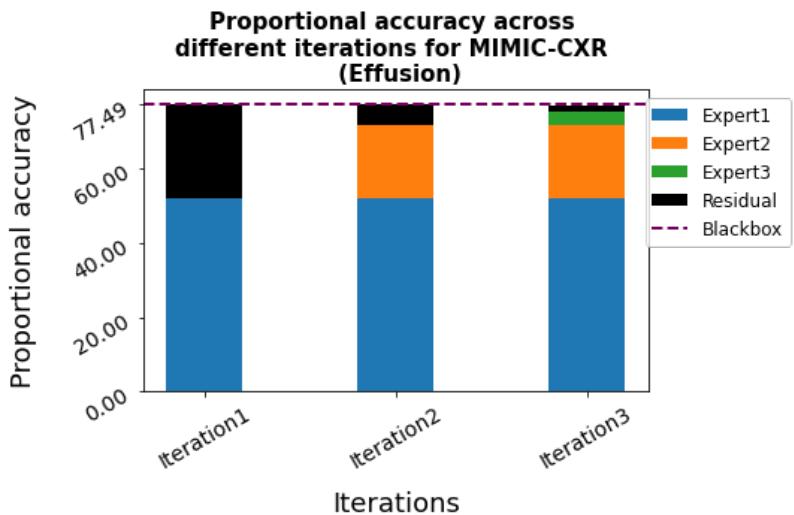
(e)

Coverages of various experts for HAM

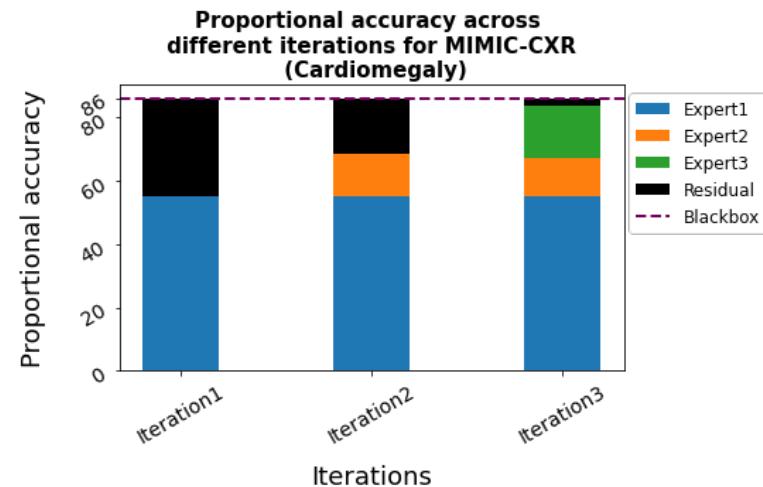
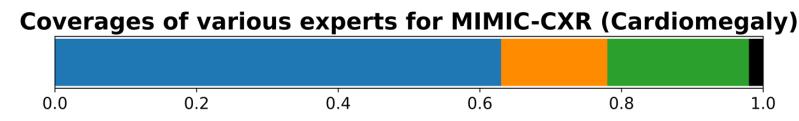
(c)

Blackbox vs Residuals across iterations for HAM10000

(f)



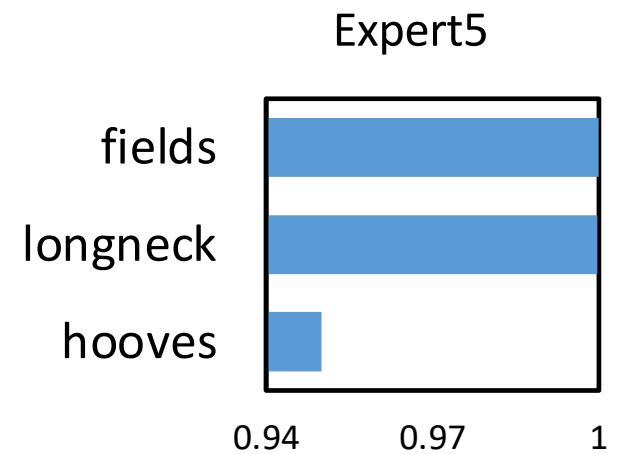
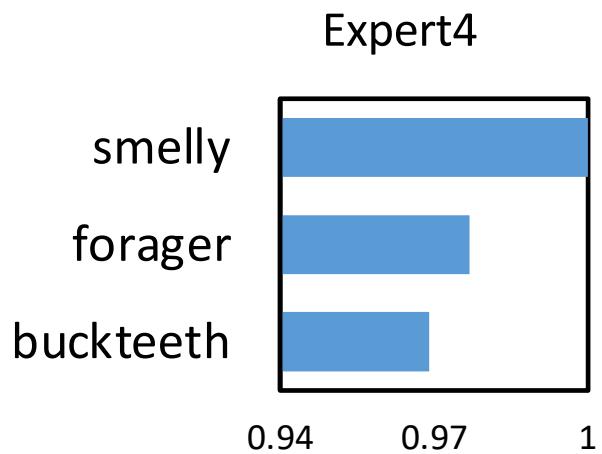
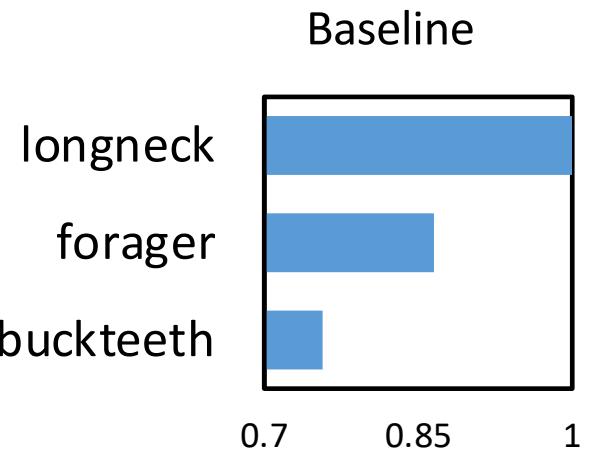
(a)



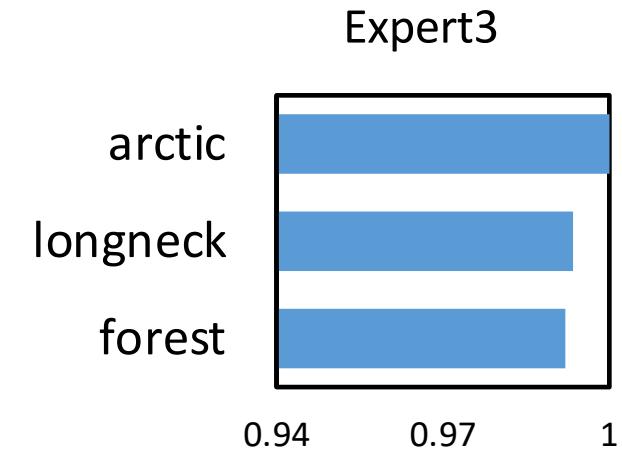
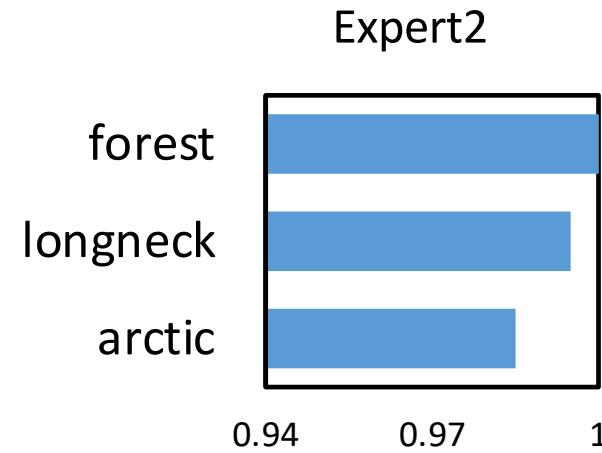
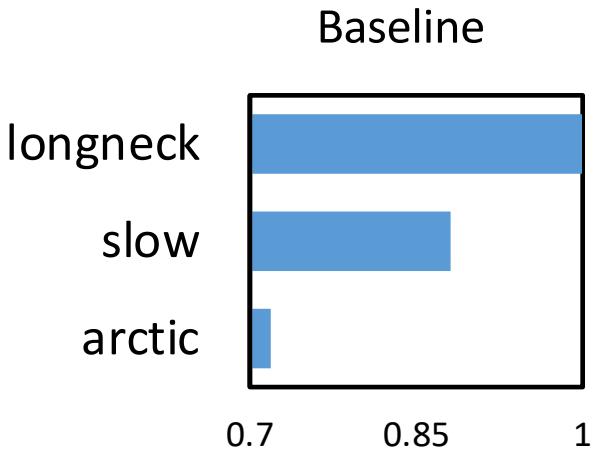
(b)

(a)

Horse



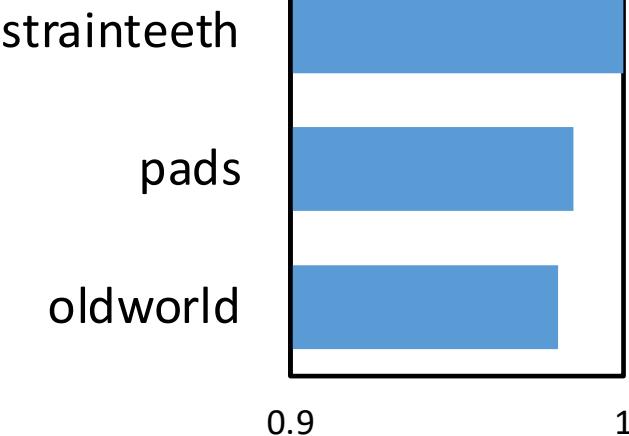
Moose



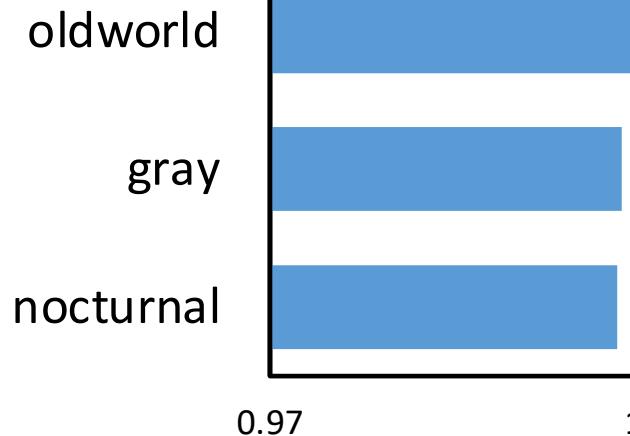
Beaver



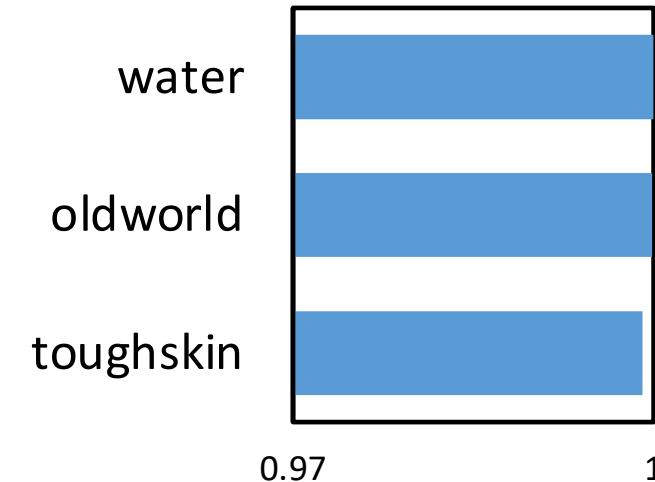
Baseline



Expert1



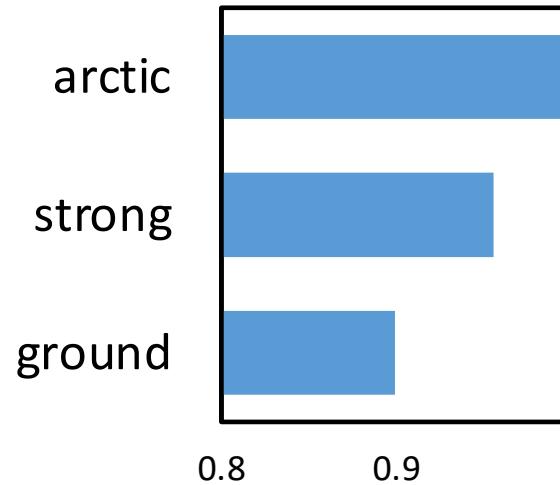
Expert4



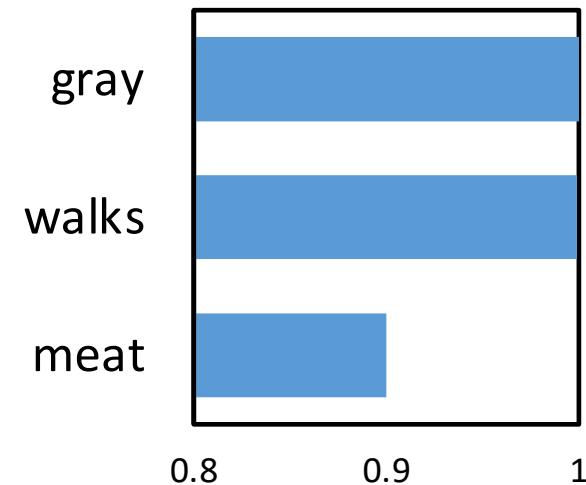
Otter



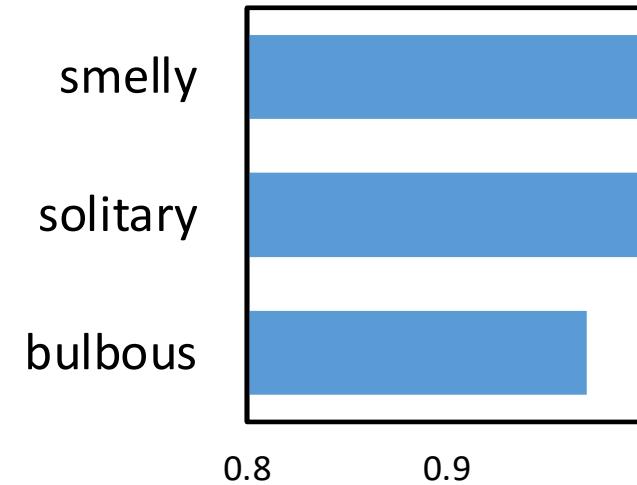
Baseline



Expert1

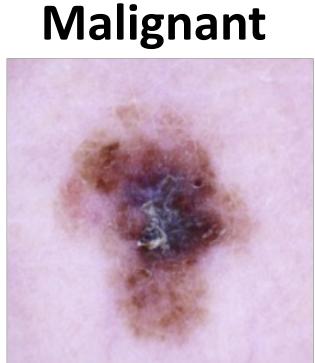
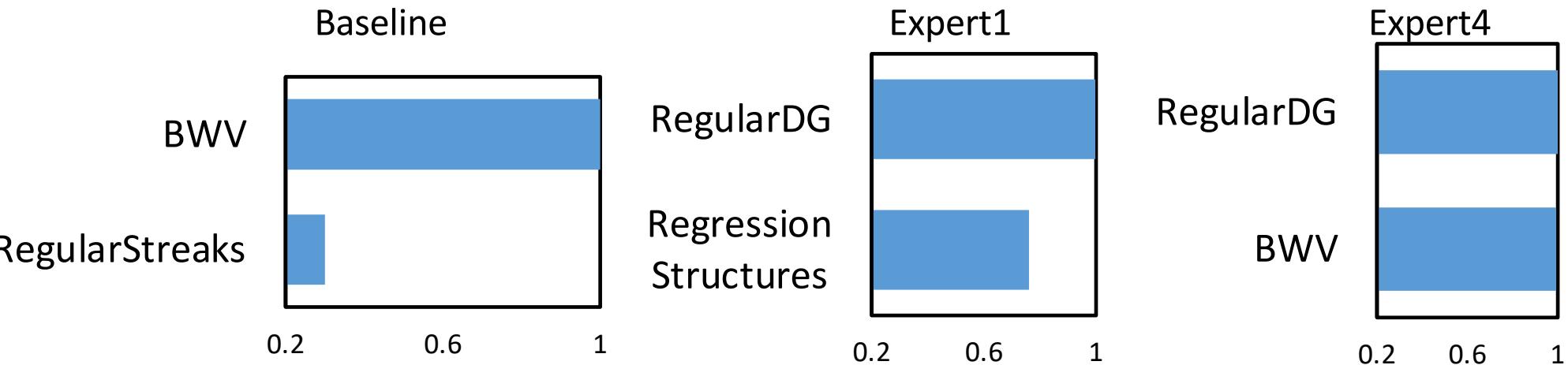


Expert2

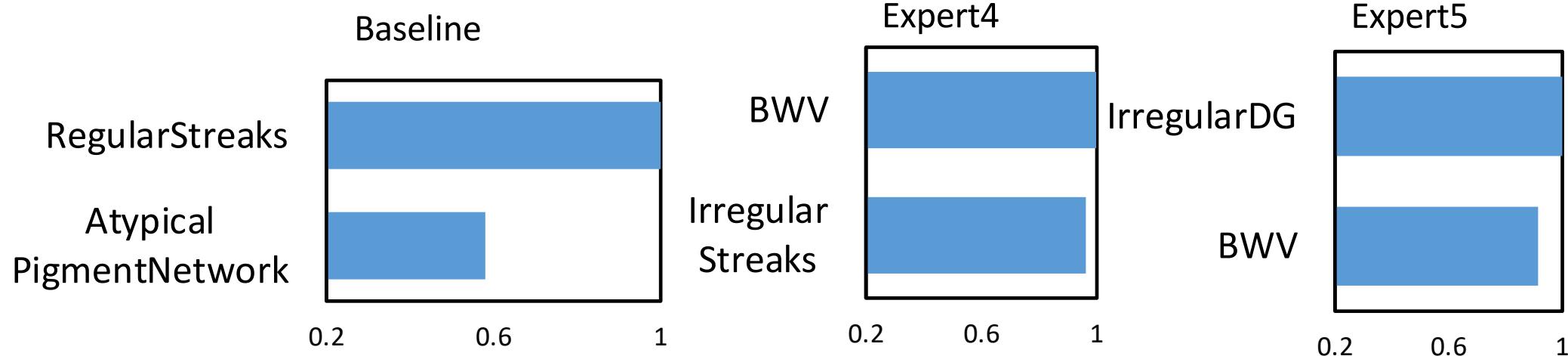




Benign



Malignant



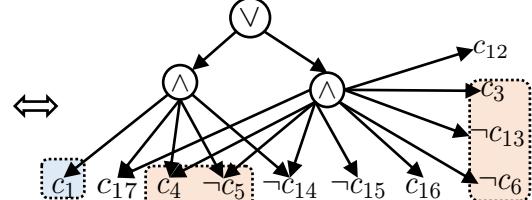
Bay breasted warbler



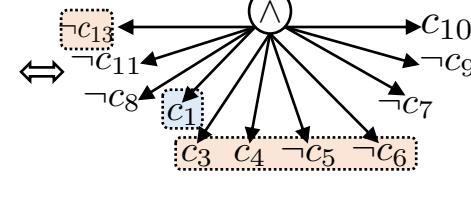
Expert1

Expert2

Expert4

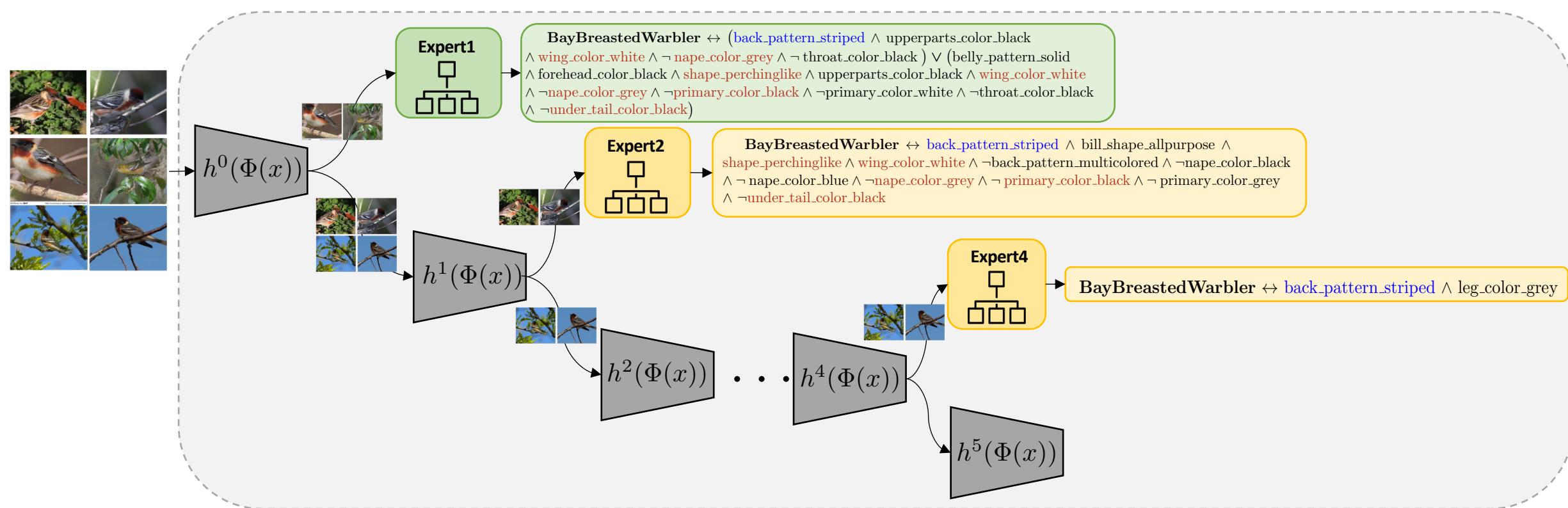


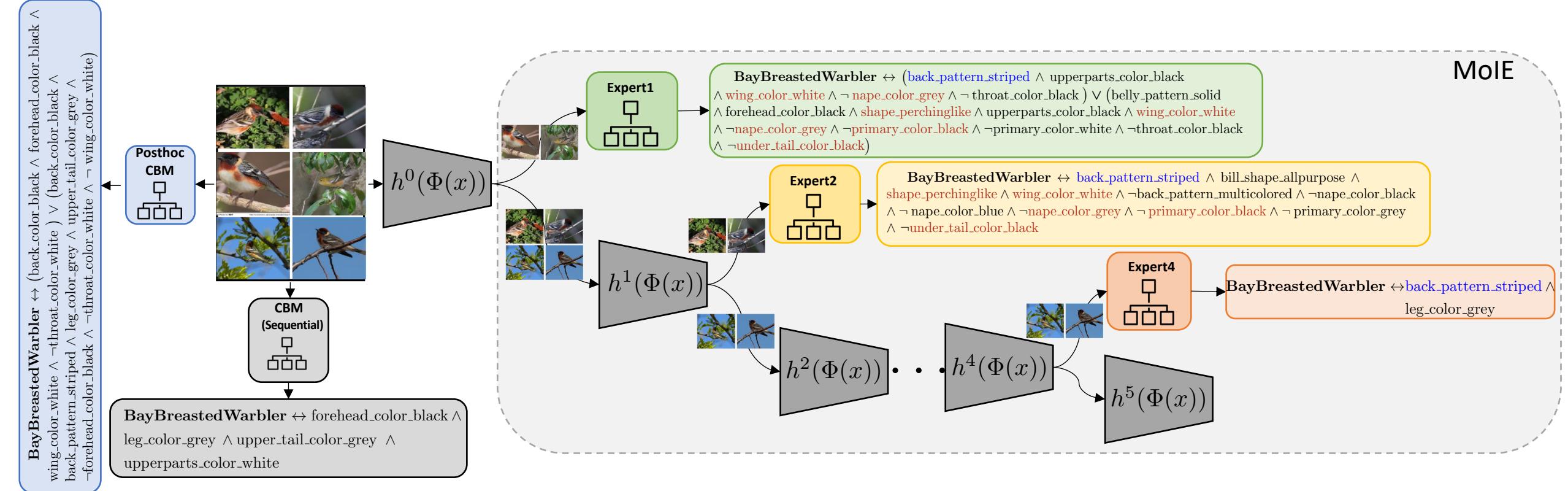
(a)

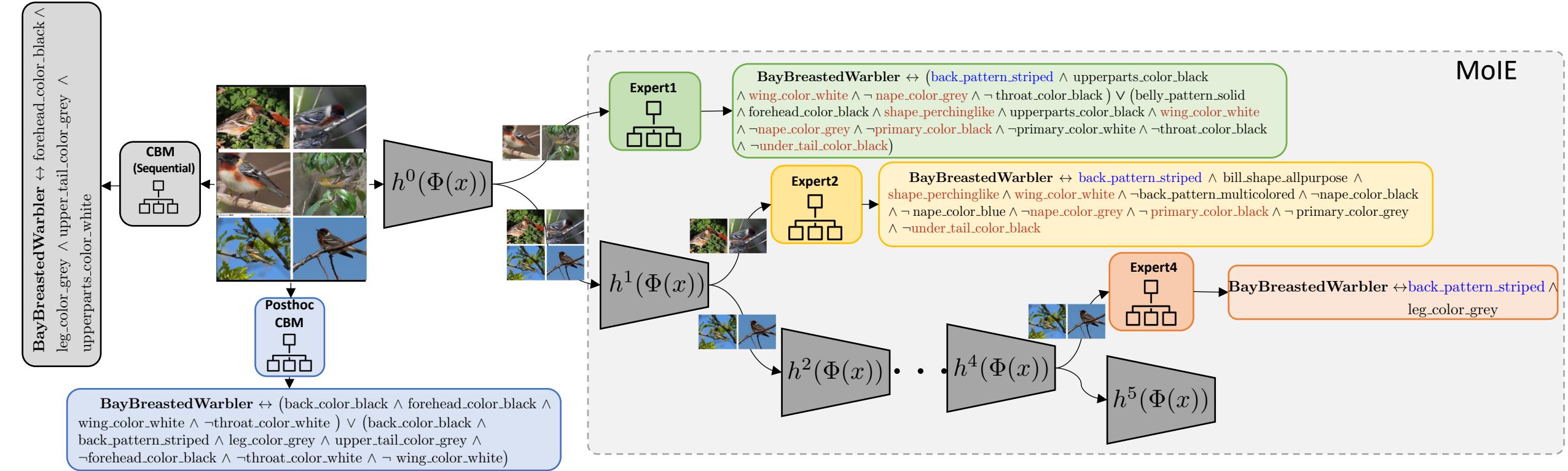


(b)

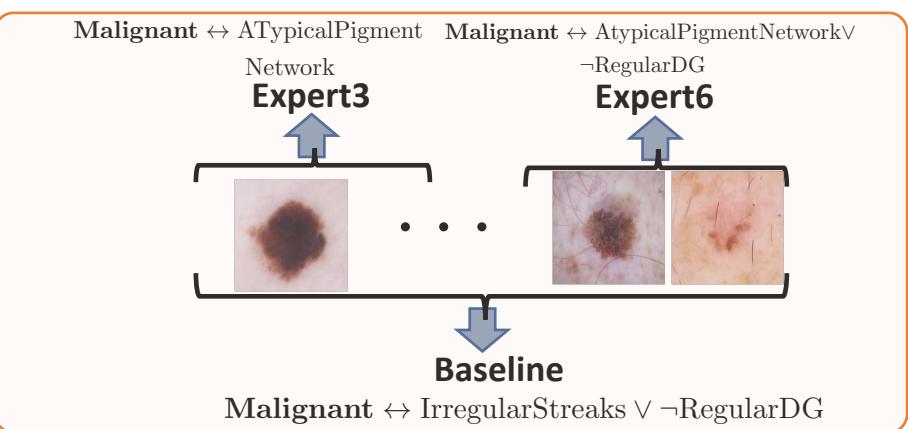
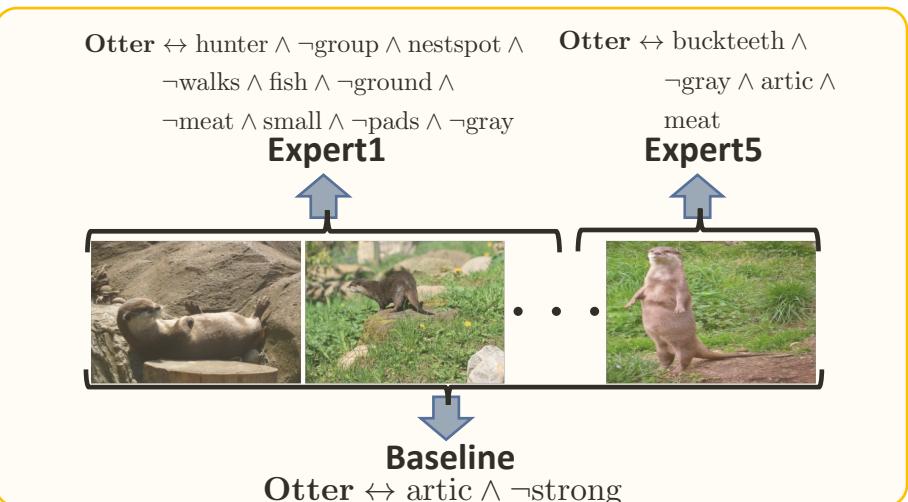
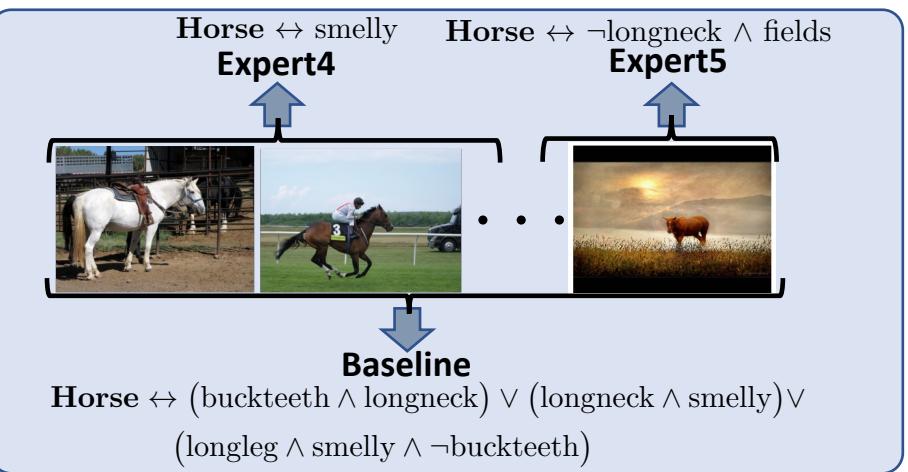
Symbol	Concept
c_1	<i>back_pattern_striped</i>
c_2	<i>leg_color_grey</i>
c_3	<i>shape_perching_like</i>
c_4	<i>wing_color_white</i>
c_5	<i>nape_color_grey</i>
c_6	<i>primary_color_black</i>
c_7	<i>primary_color_grey</i>
c_8	<i>back_pattern_multicolored</i>
c_9	<i>nape_color_black</i>
c_{10}	<i>bill_shape_allpurpose</i>
c_{11}	<i>nape_color_blue</i>
c_{12}	<i>forehead_color_black</i>
c_{13}	<i>undertail_color_black</i>
c_{14}	<i>throat_color_black</i>
c_{15}	<i>primary_color_white</i>
c_{16}	<i>belly_pattern_solid</i>
c_{17}	<i>upperparts_color_black</i>



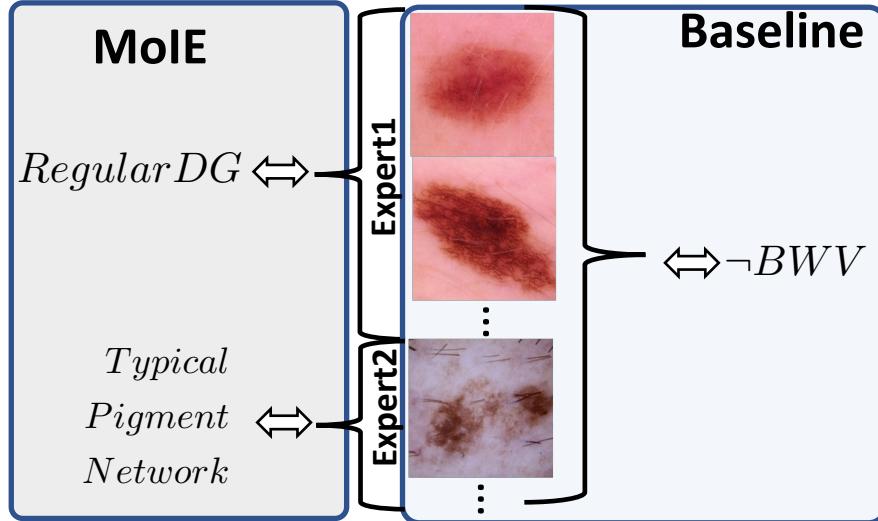




Local Logical explanations. Generic by the Baseline. Not clear which concept for which sample

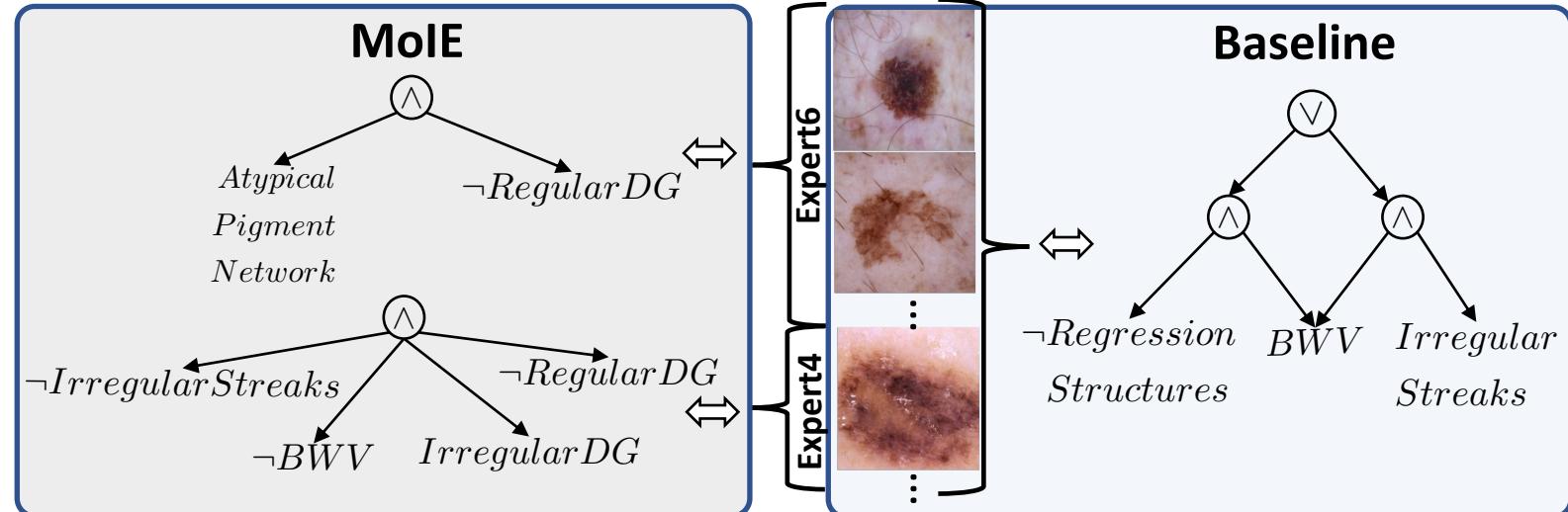


Class: Benign

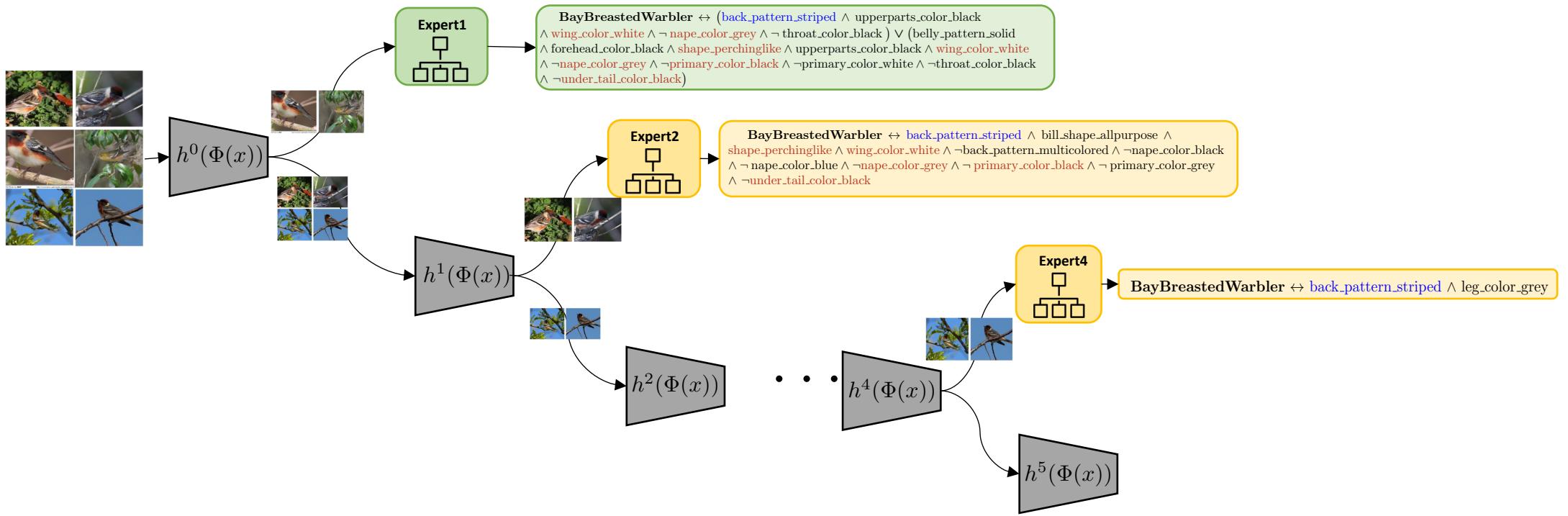


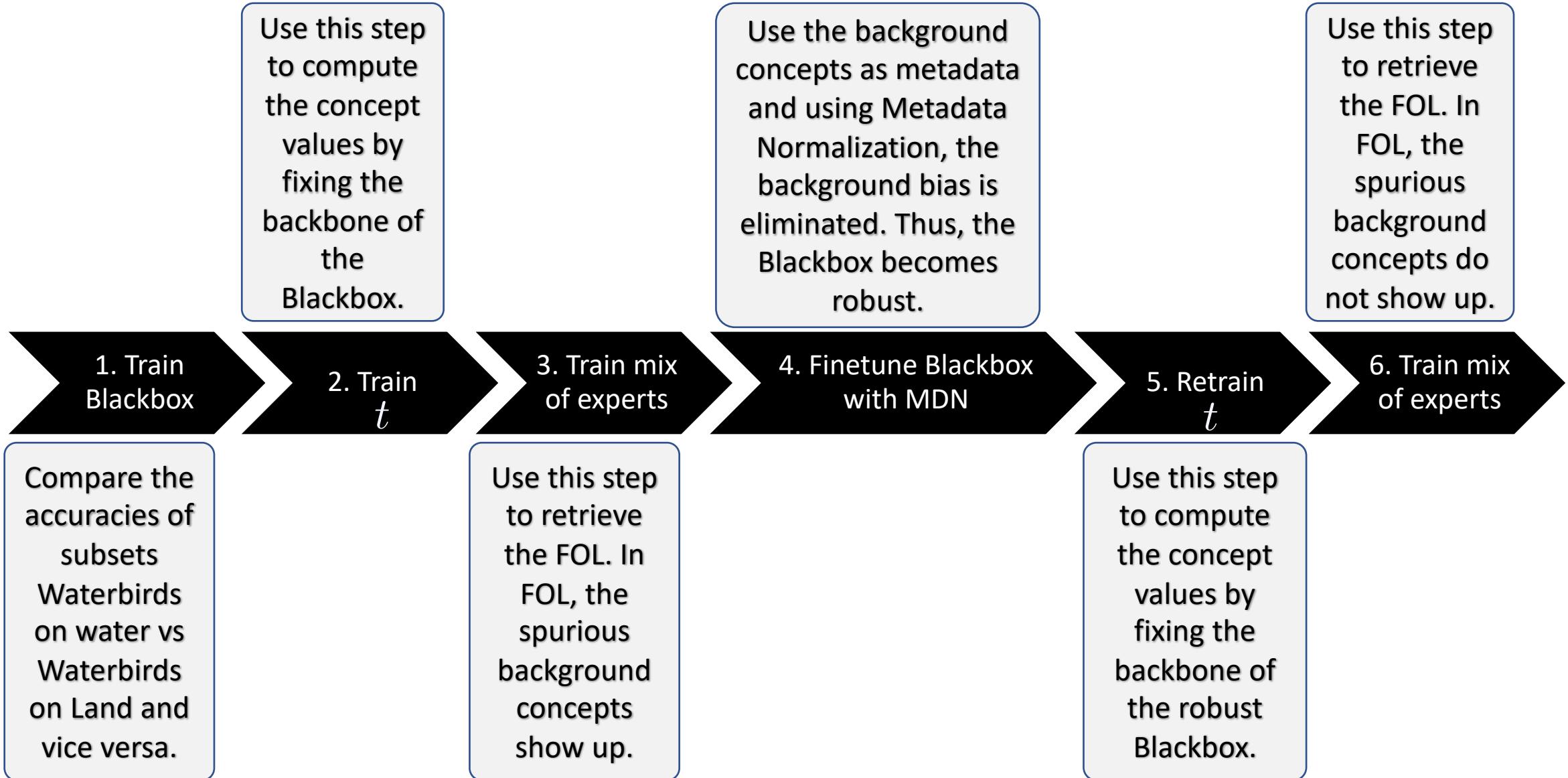
(c)

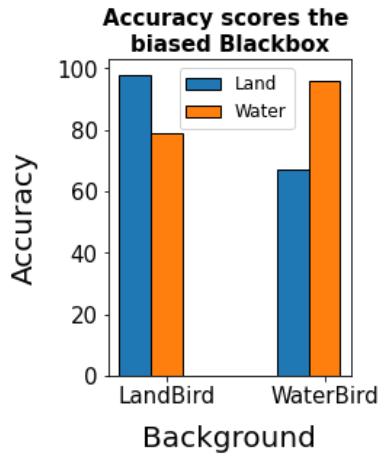
Class: Malignant



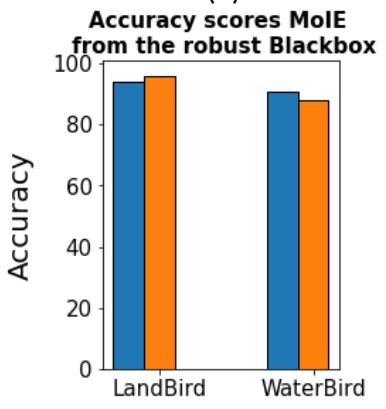
(d)



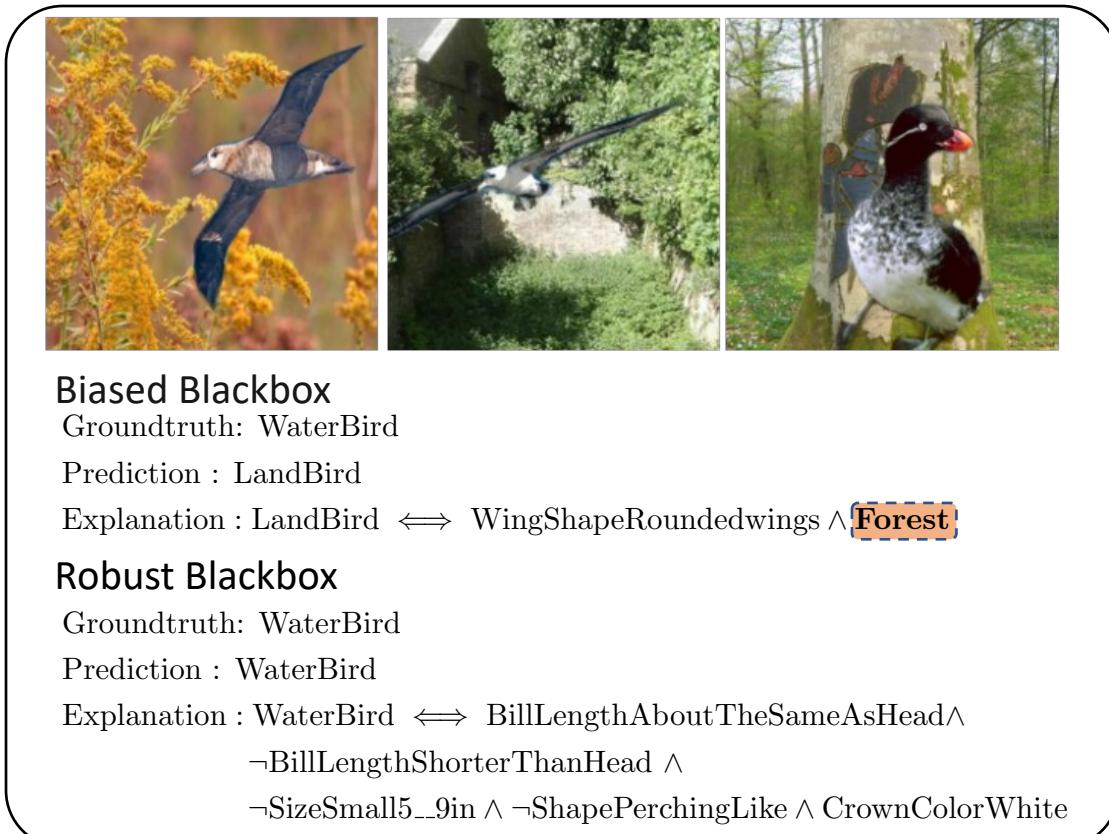




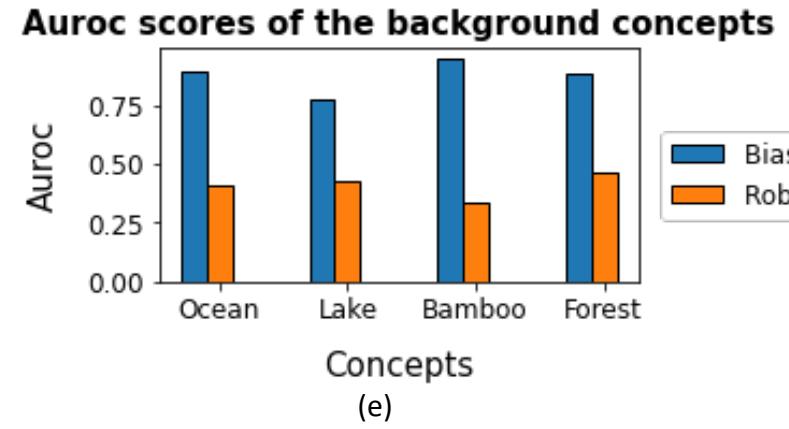
(a)



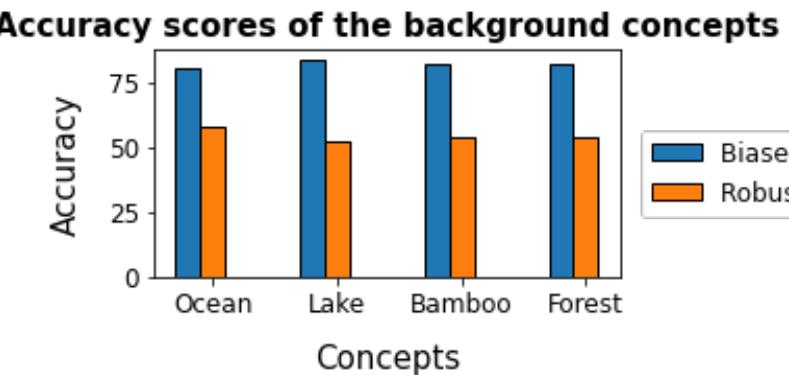
(b)



(c)

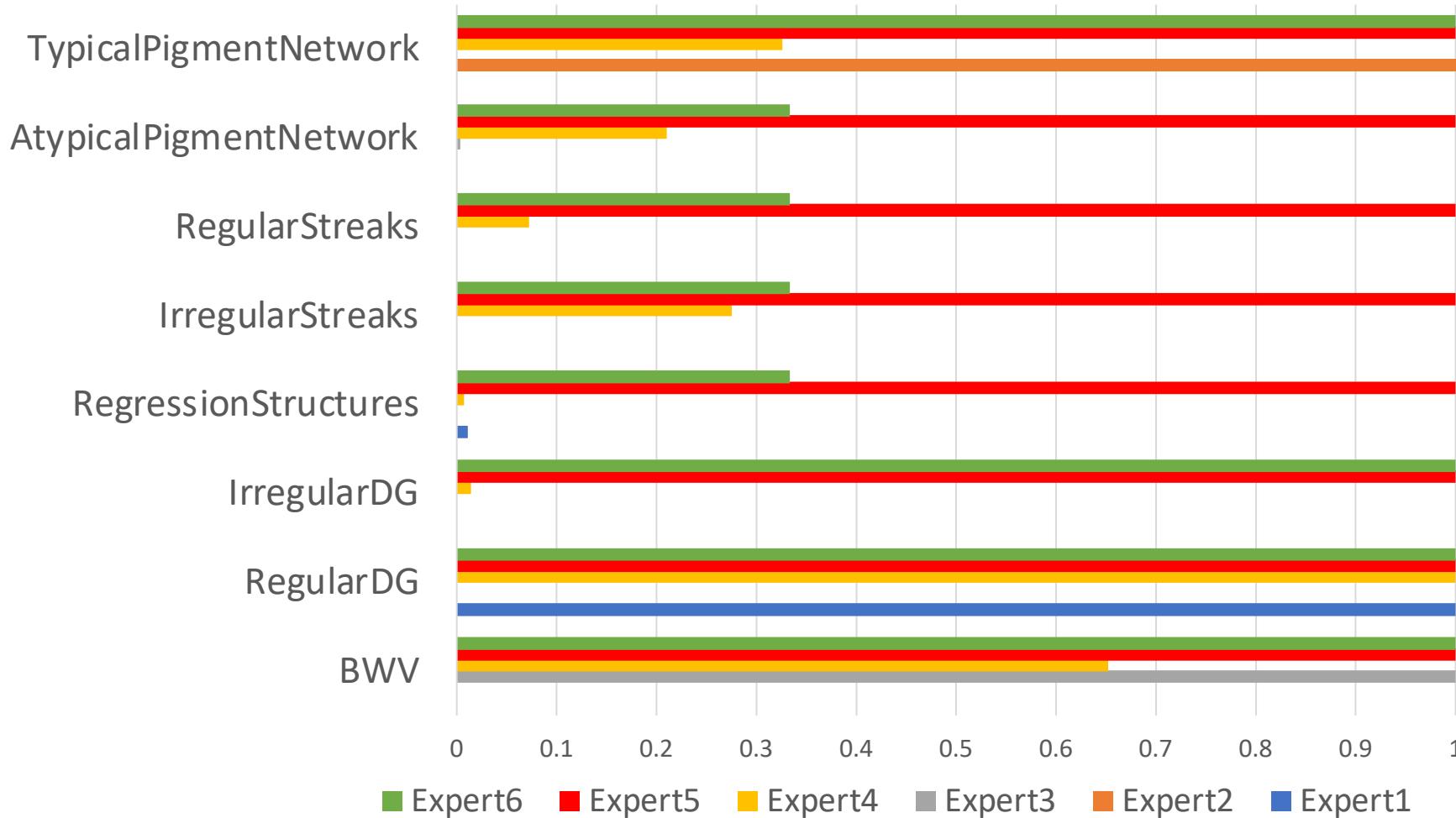


(e)

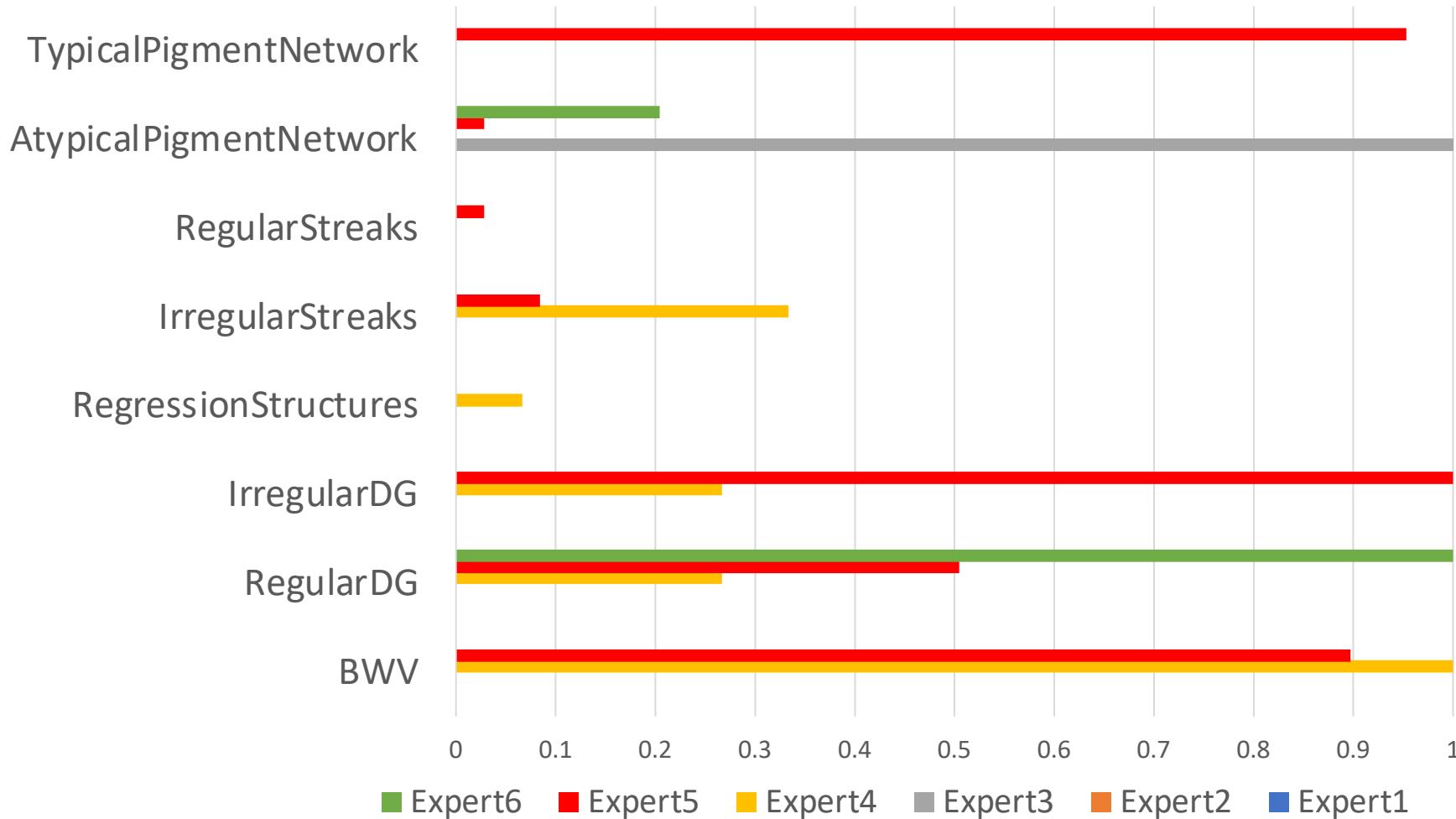


(d)

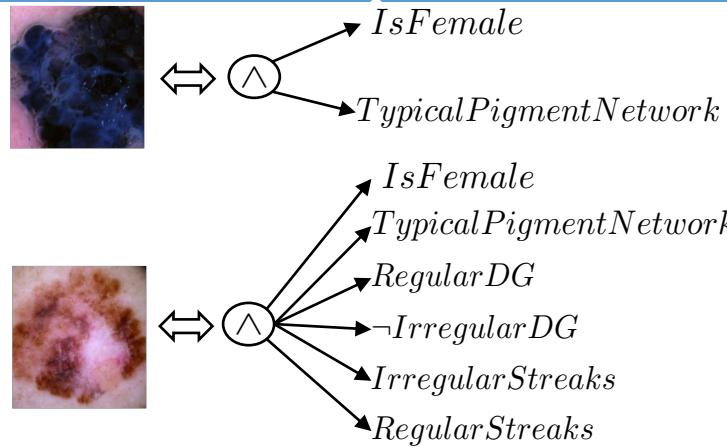
Benign



Malignant

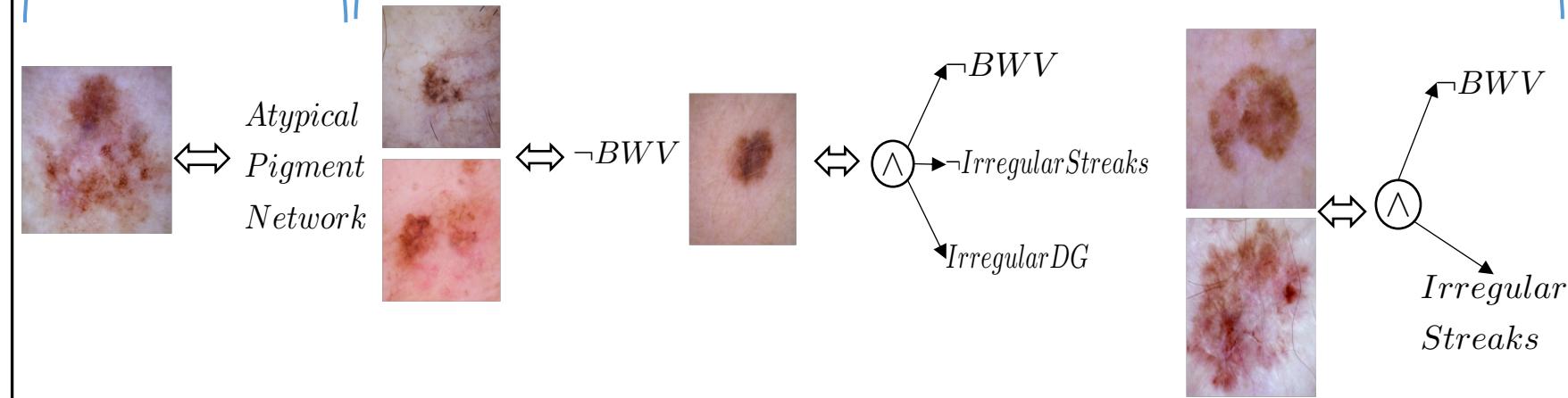


Baseline
(Interpretable by design)



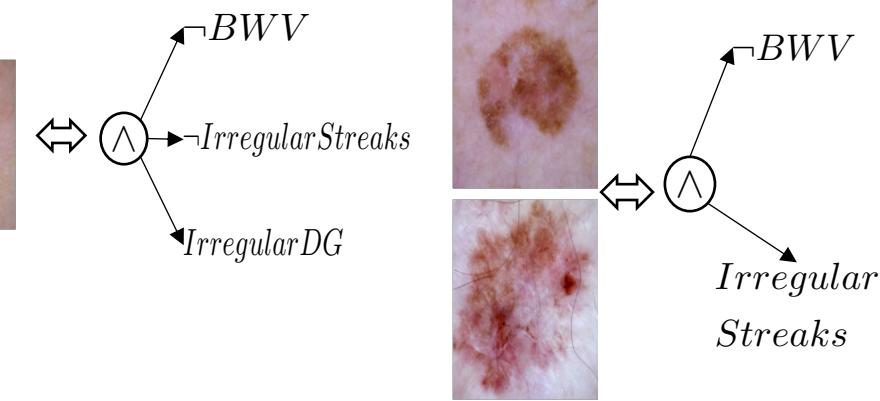
(a)

Expert2

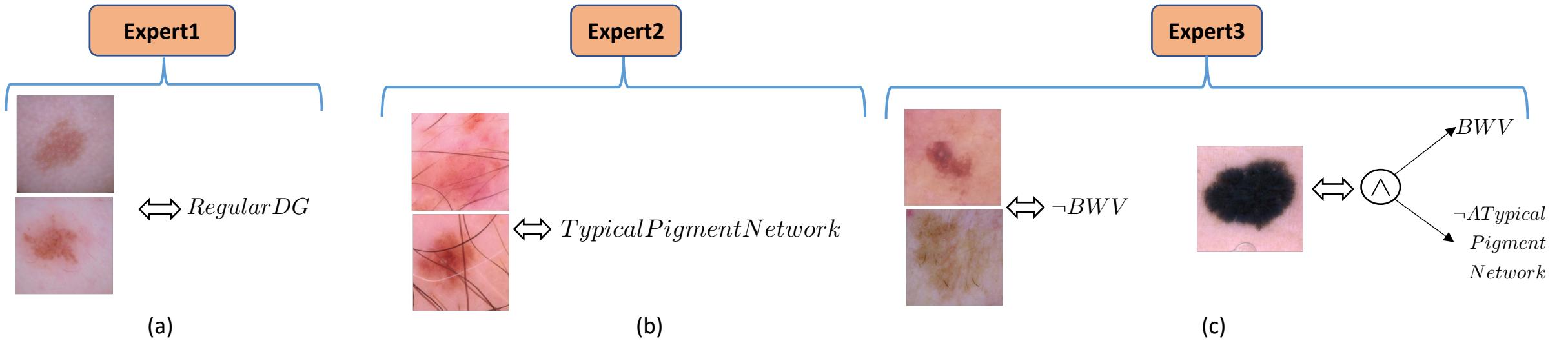


(b)

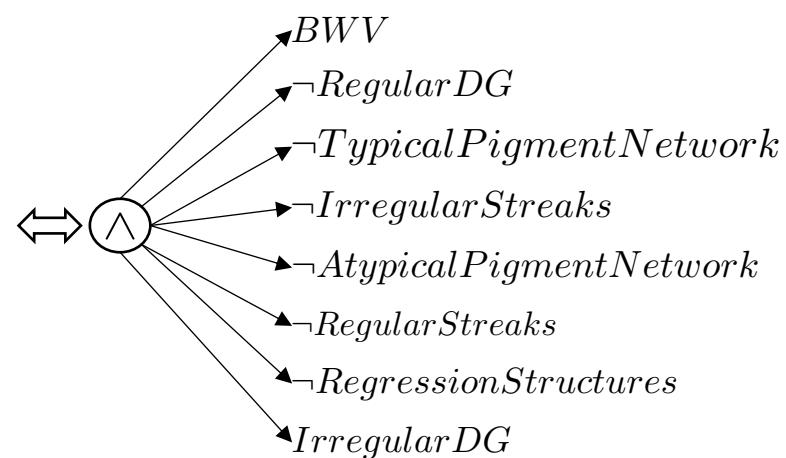
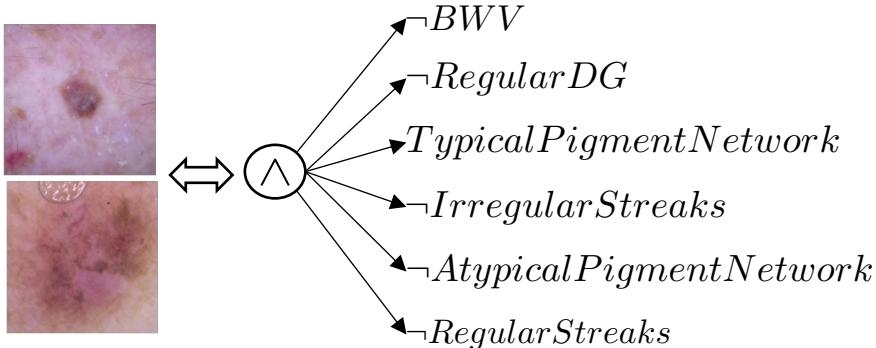
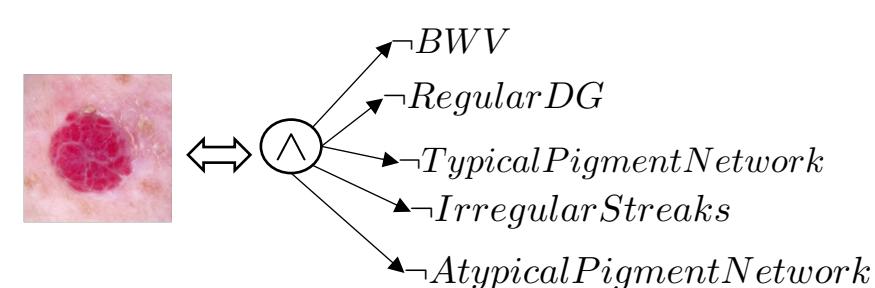
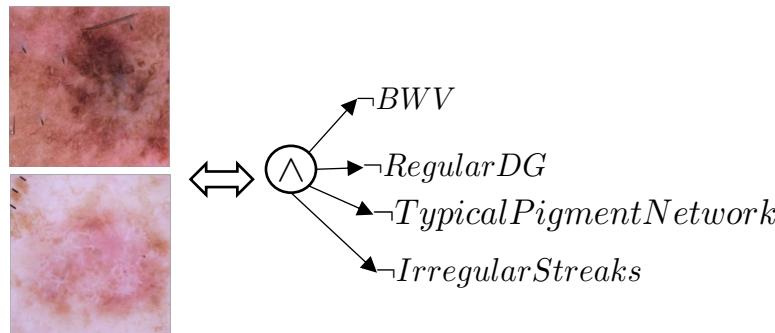
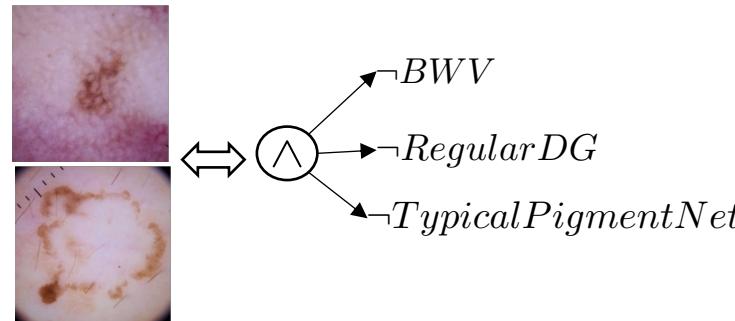
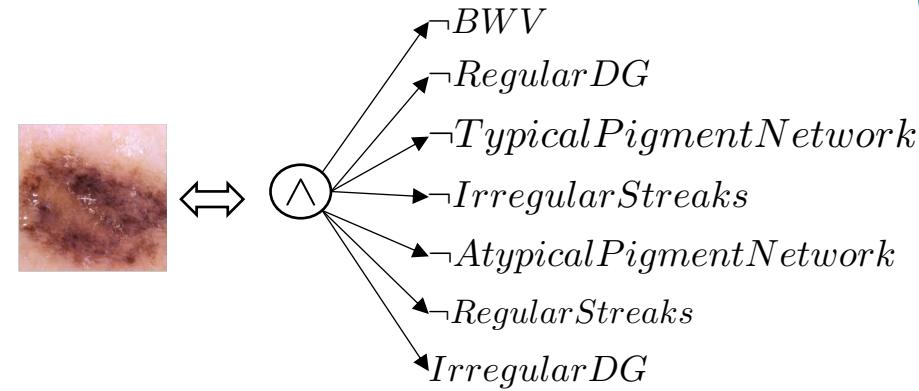
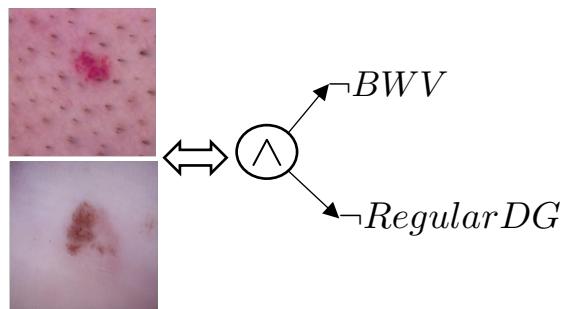
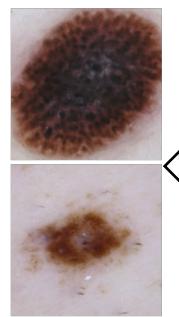
Expert4



(c)



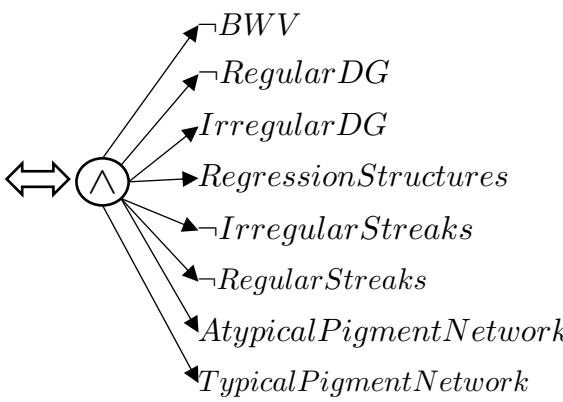
Expert4



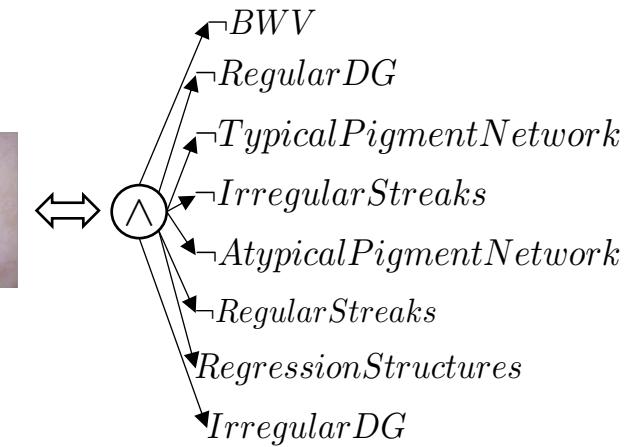
(d)

Expert6

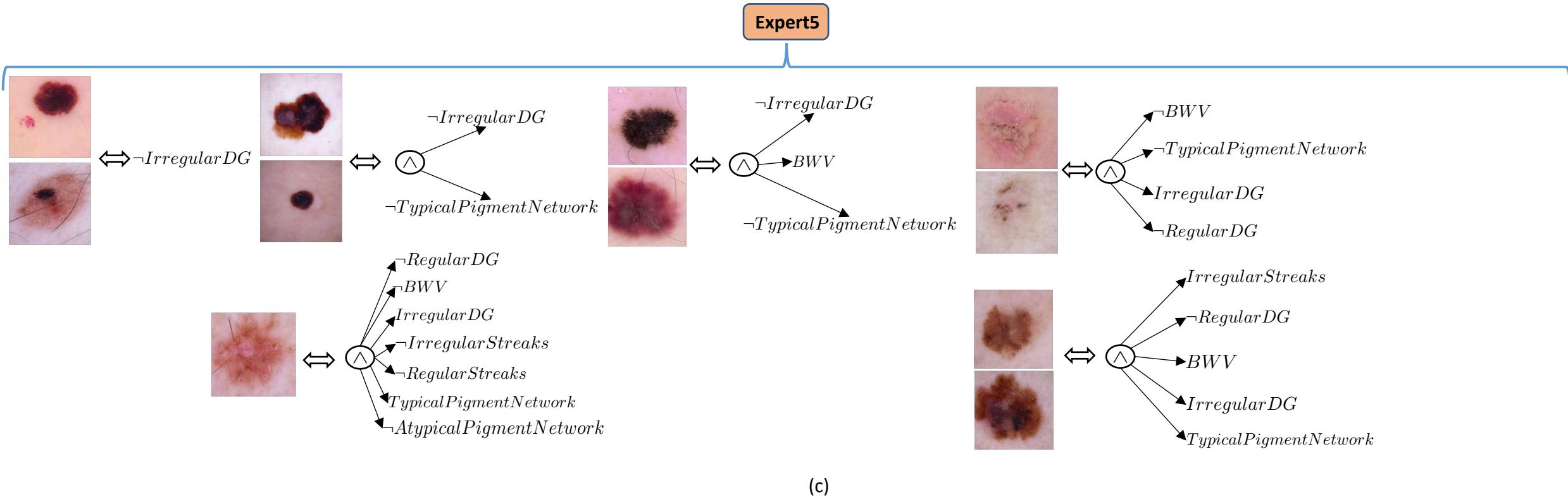
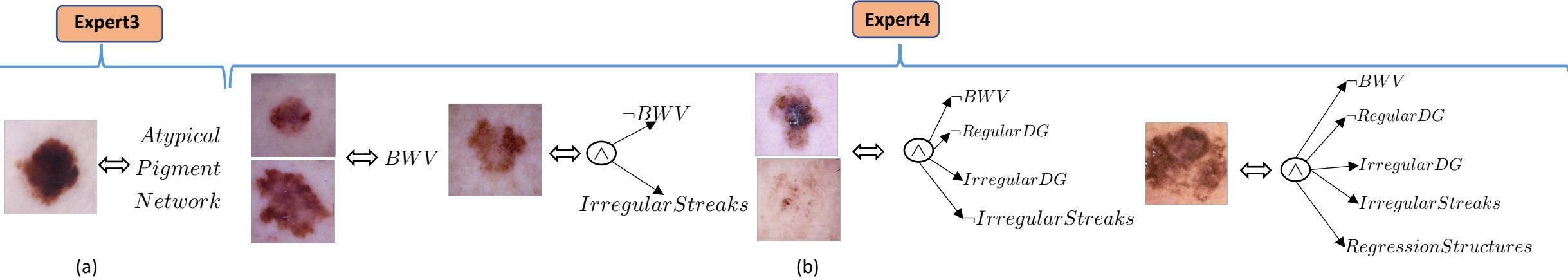
Expert5

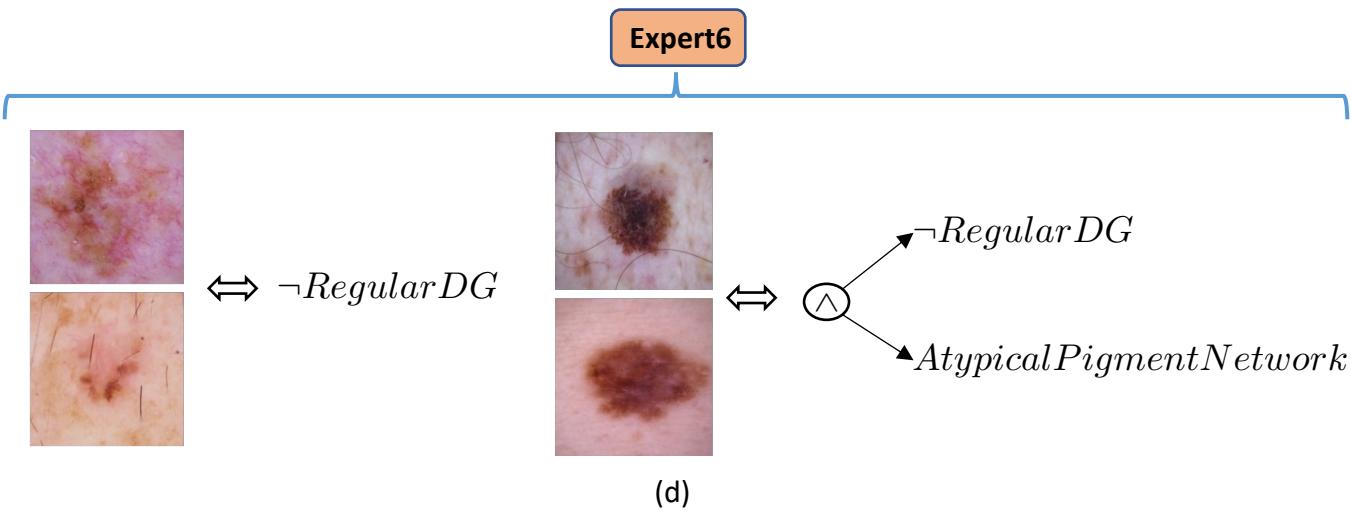


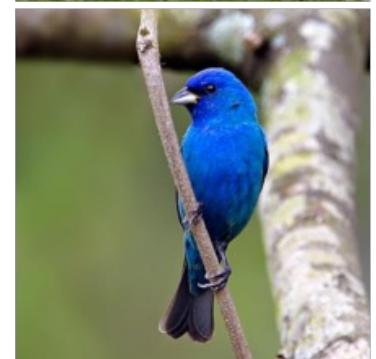
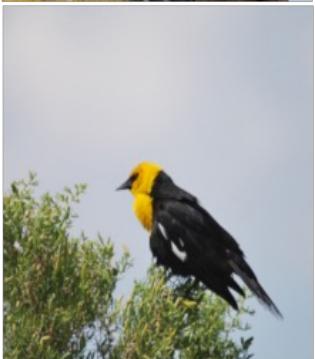
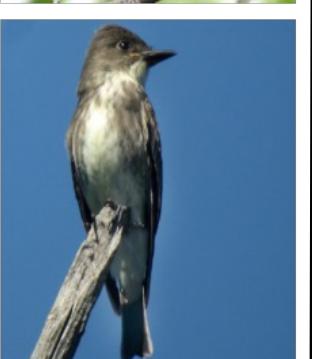
(e)



(f)



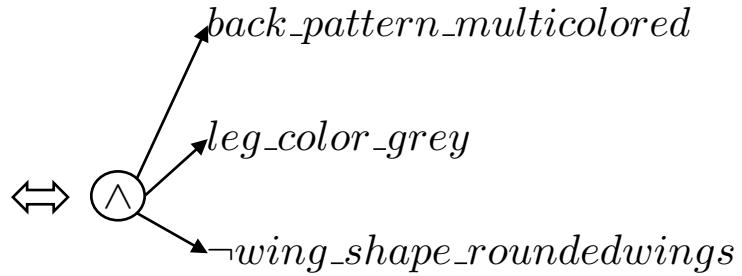


Expert1**Expert2****Expert3****Expert4****Expert5****Expert6****Final Residual**

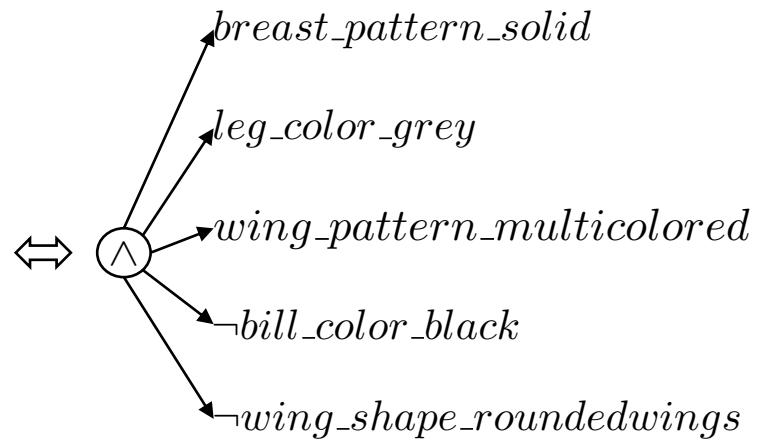


Painted Bunting

Expert1

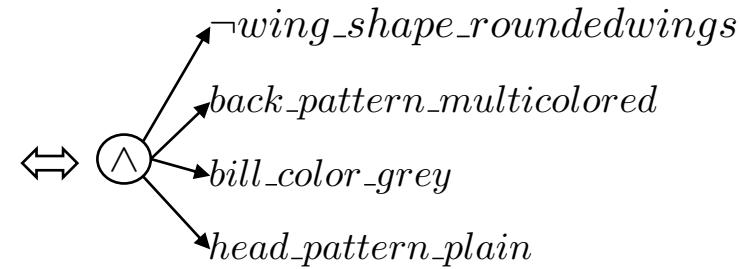


Expert2

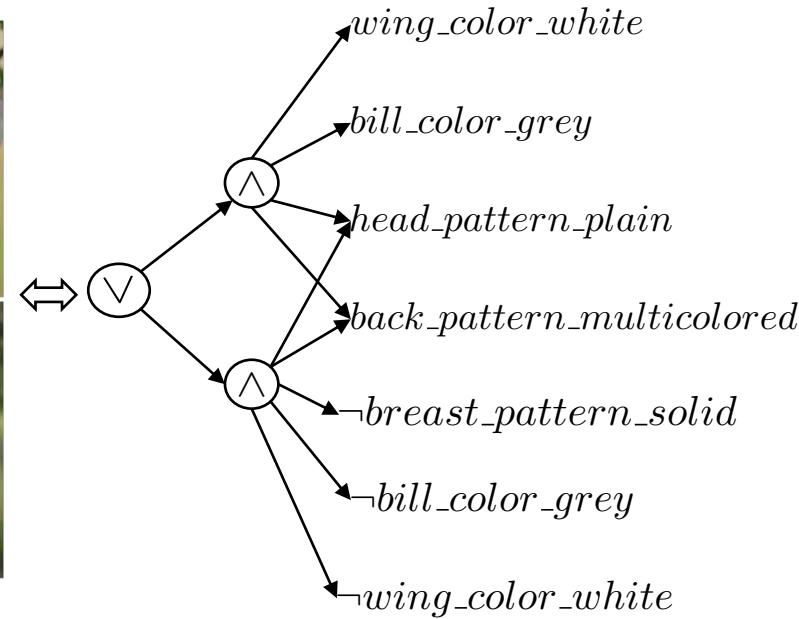


Baltimore Oriole

Expert1

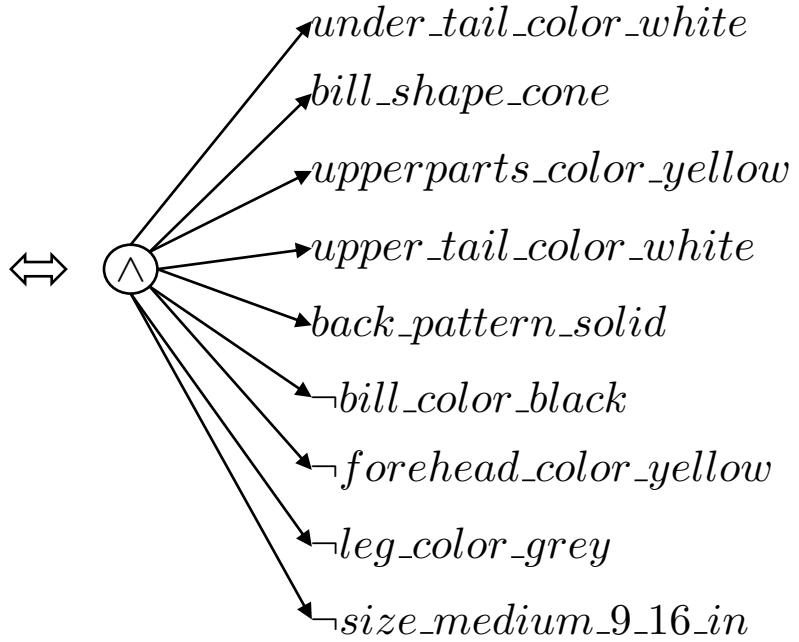


Expert3

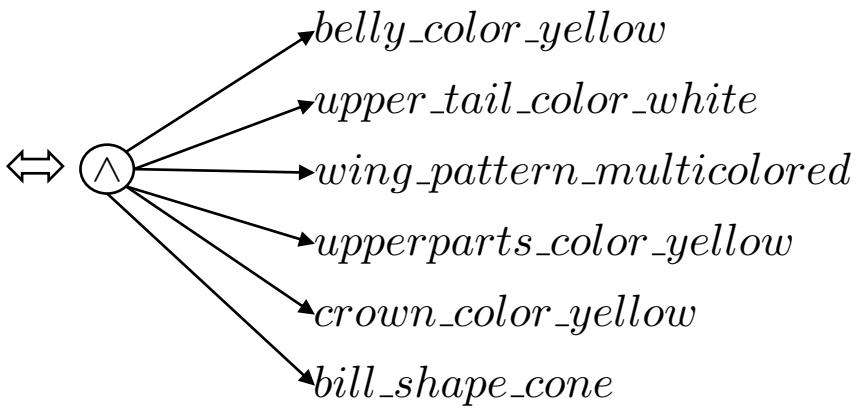


American Goldfinch

Expert1



Expert2



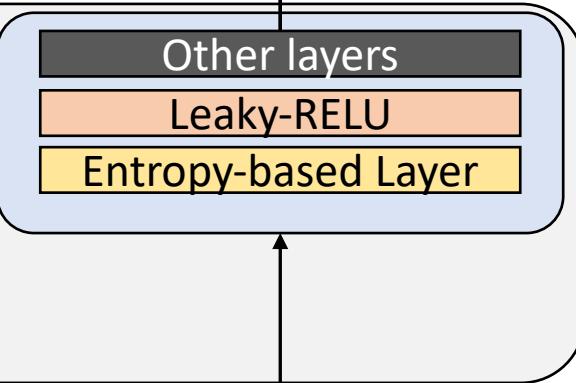
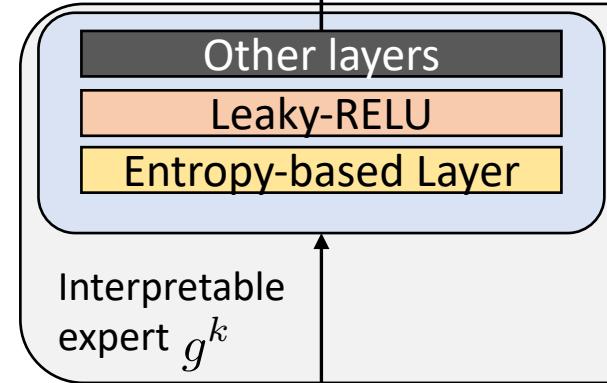
Predict 1st class

FOL 1st class

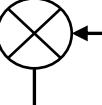
Predict rth class

FOL rth class

Blackbox
Logits



Interpretable
expert g^k

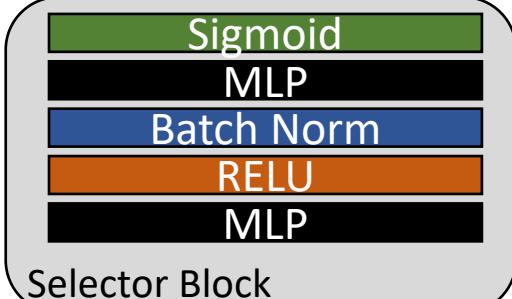


$$\pi^k > 0.5$$



$$\pi^k > 0.5$$

Input Concepts



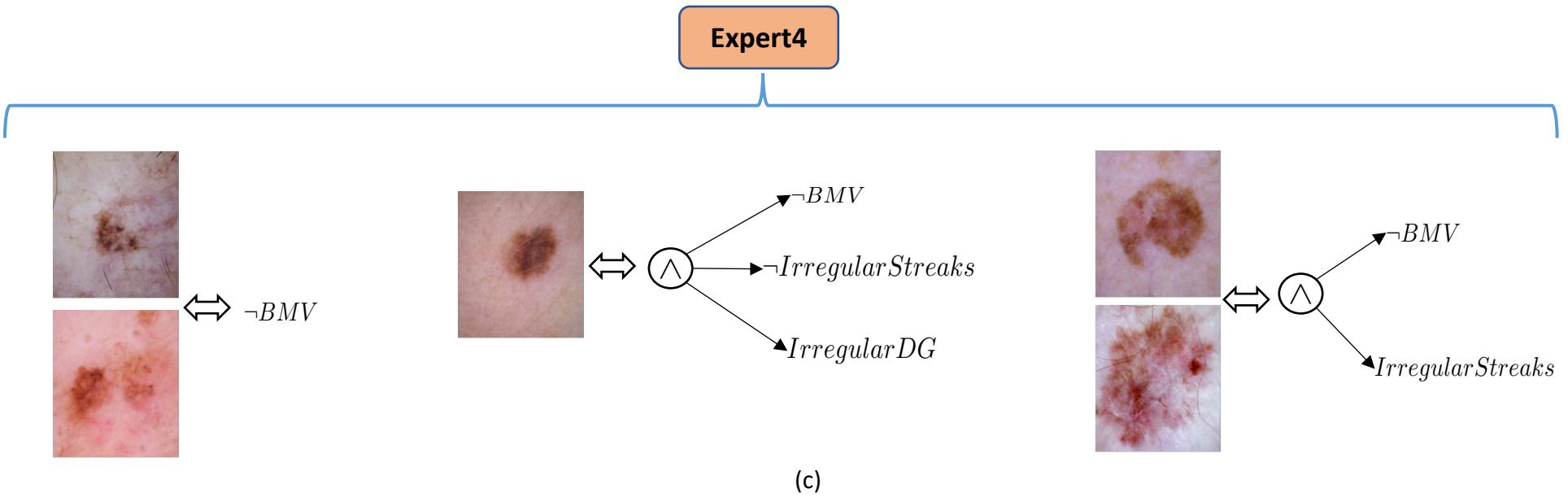
Blackbox for
iteration $k+1$

h^k

Feature
Backbone
 ϕ

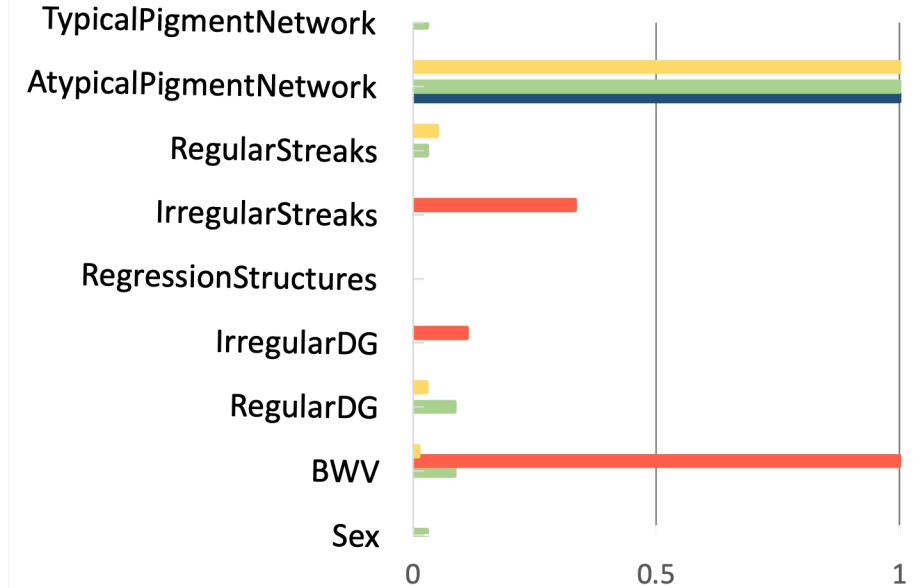
f^k

$\pi^k \leq 0.5$



MALIGNANT

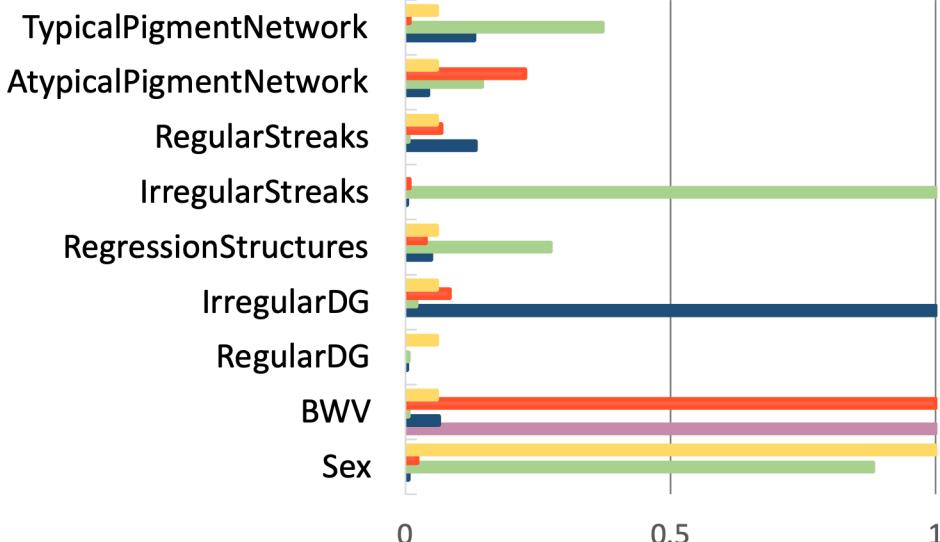
Expert5 Expert4 Expert3 Expert2 Expert1

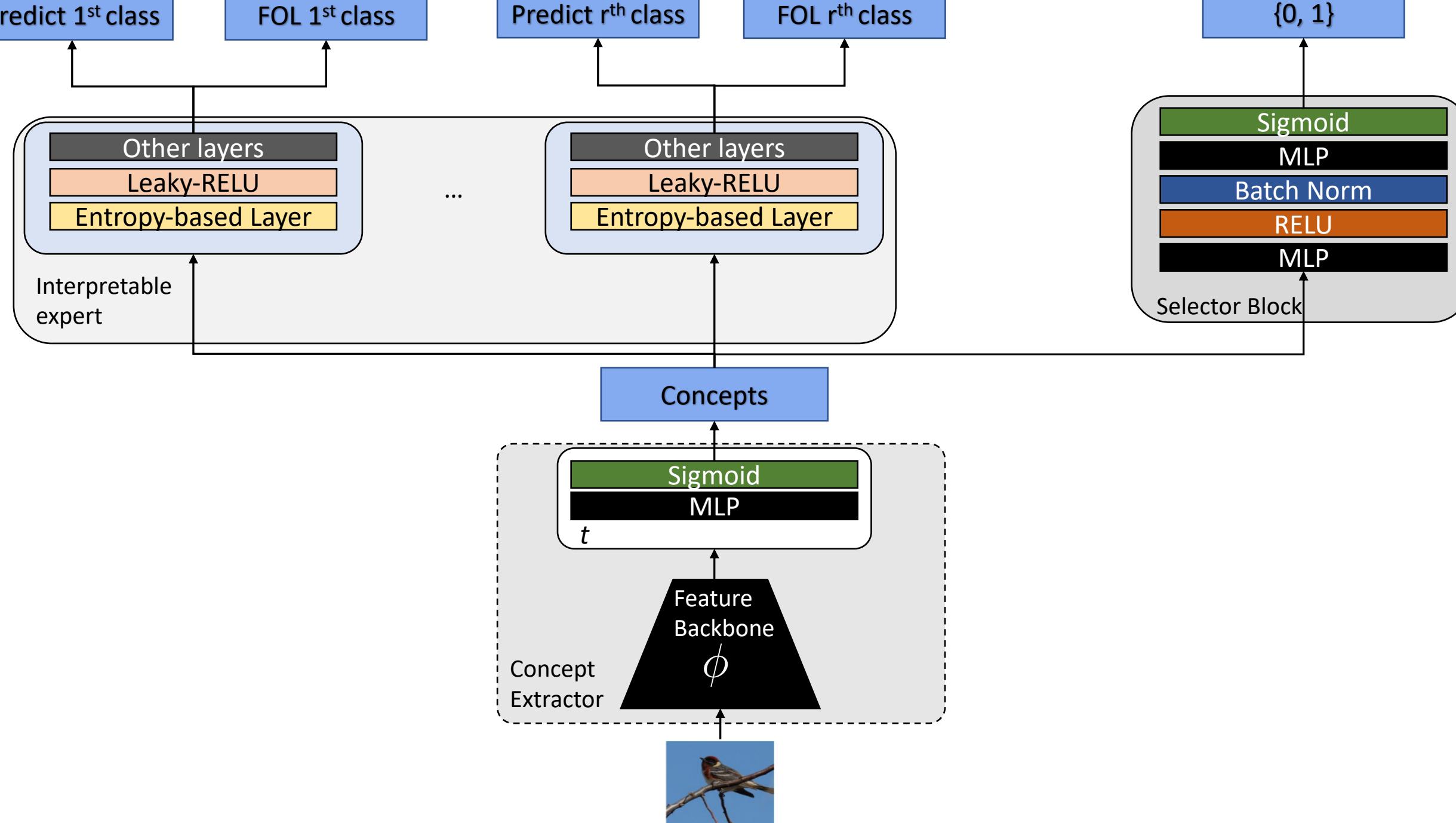


(a)

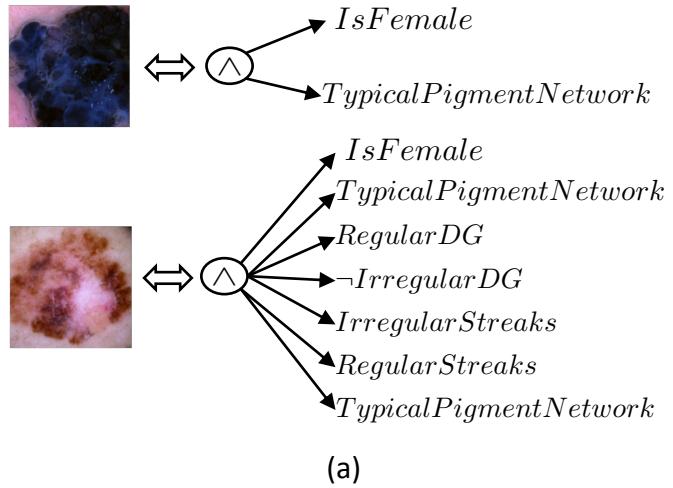
BENIGN

Expert5 Expert4 Expert3 Expert2 Expert1

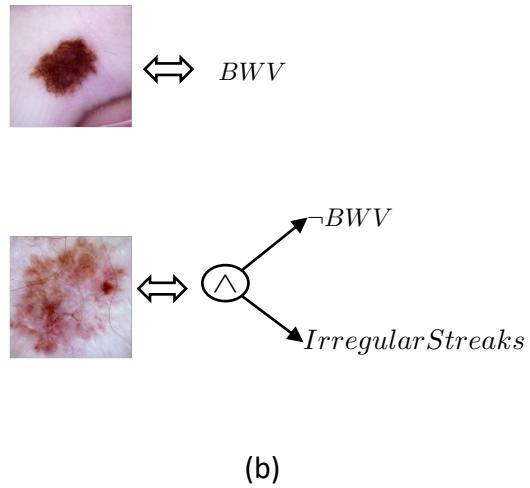




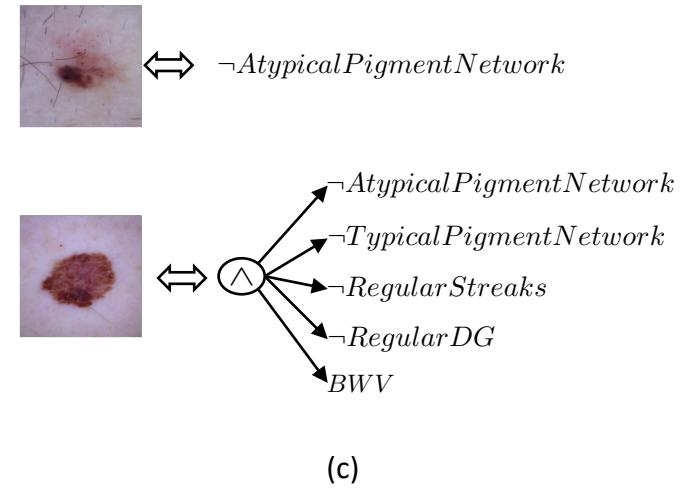
Interpretable by design

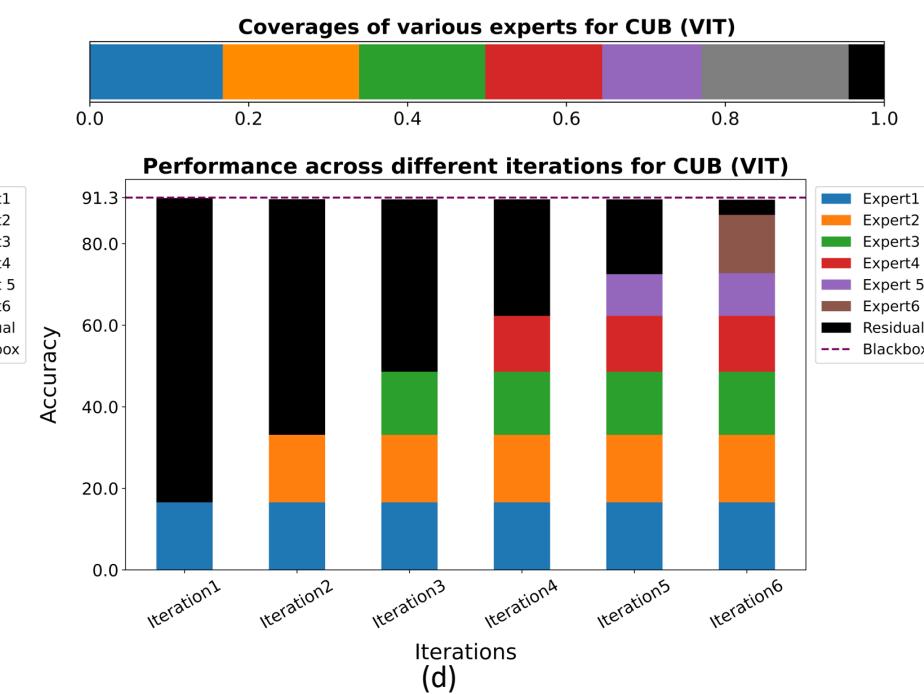
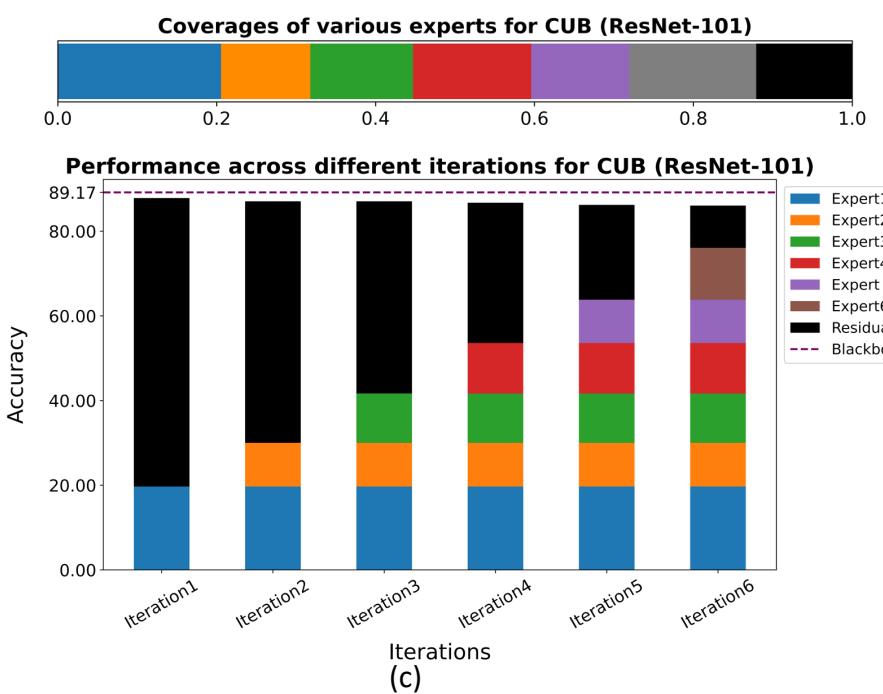
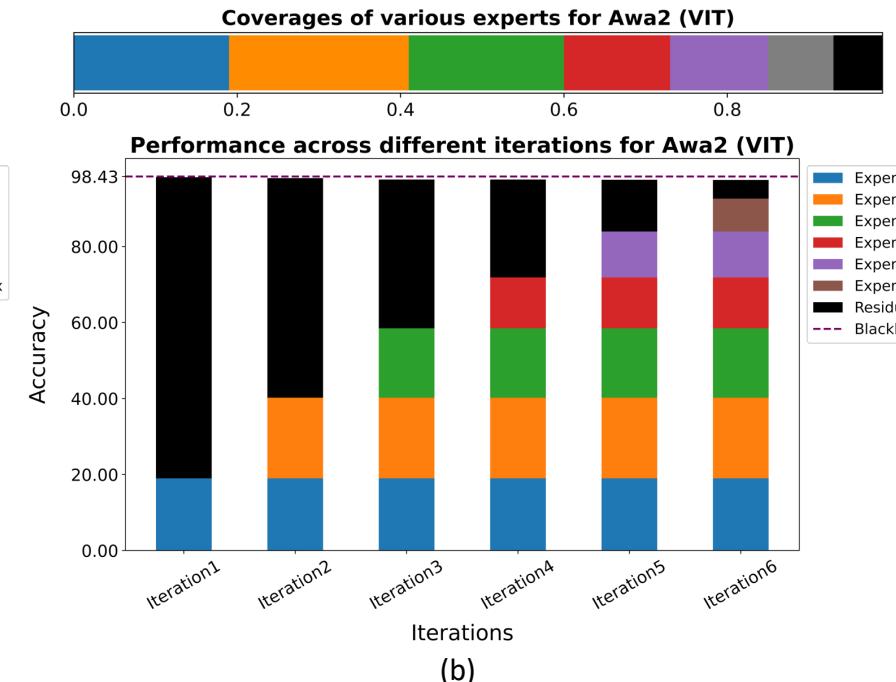
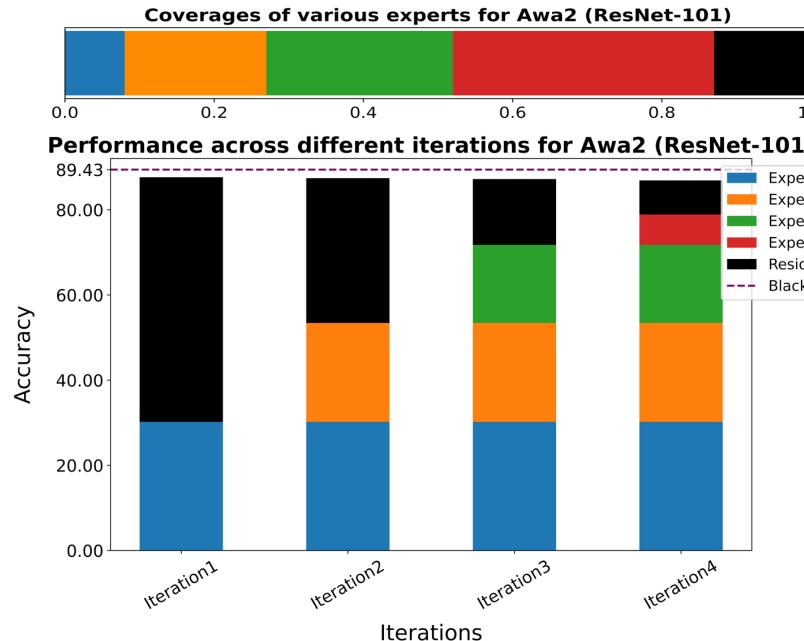


Expert4

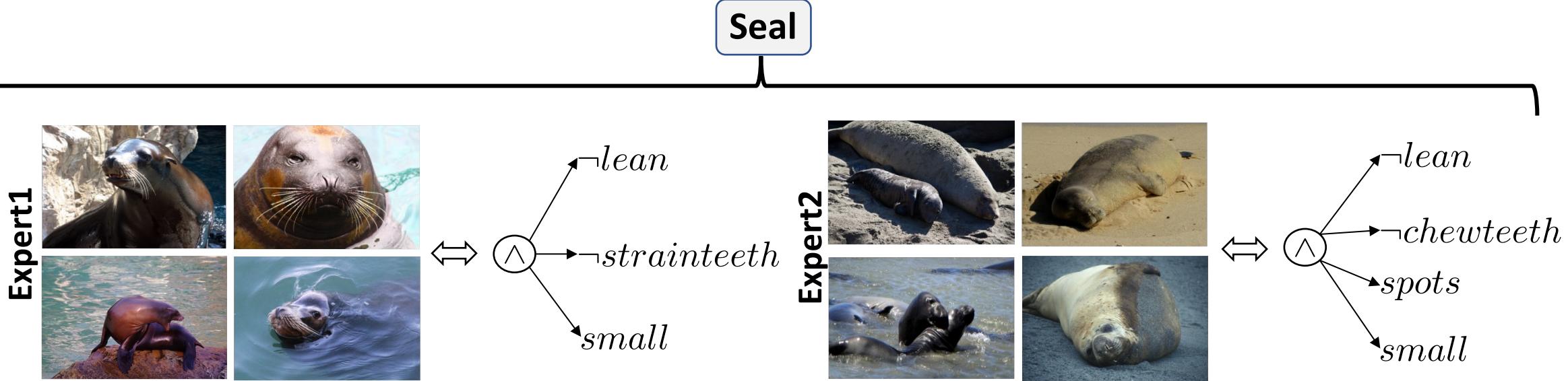


Expert5

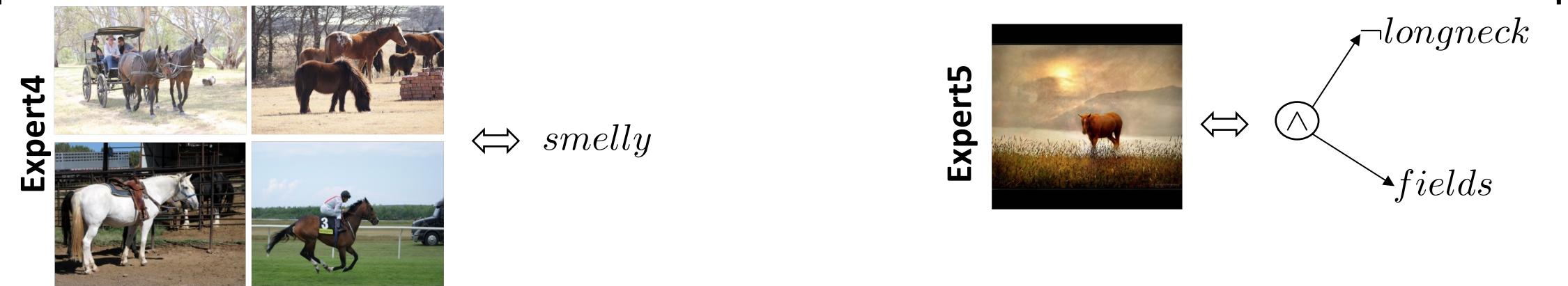




Seal



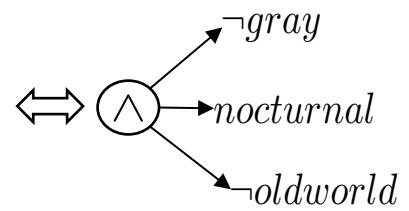
Horse



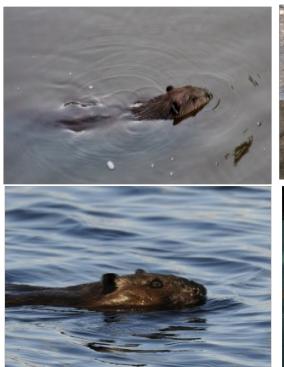
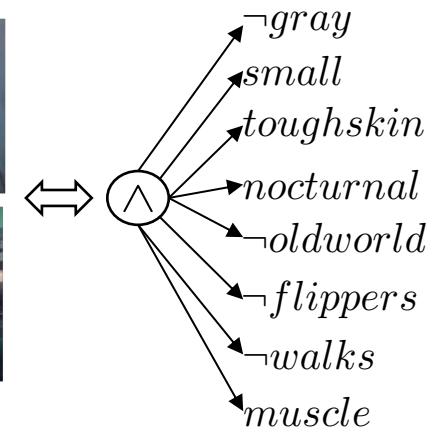
Beaver



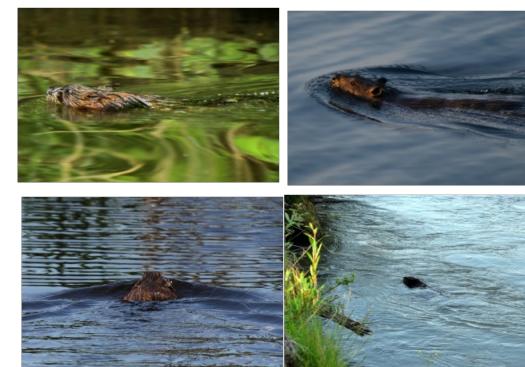
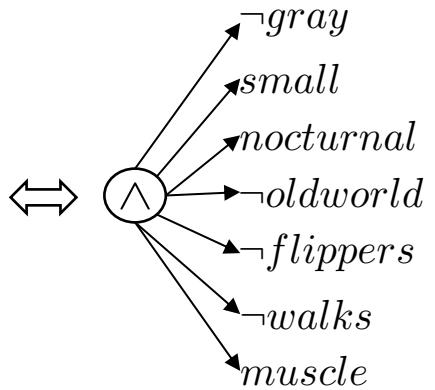
Expert1



Expert1



Expert1

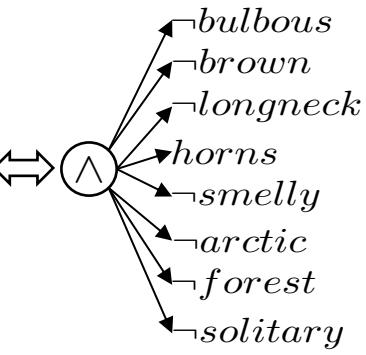


Expert4

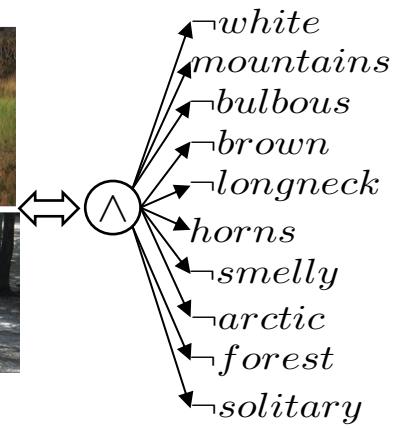
↔ *water*

Antelope

Expert2



Expert2



Expert3

